

TAU Performance System®

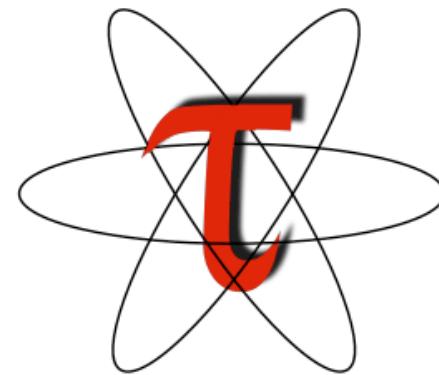
Leap to Petascale Workshop, ALCF

Room A, ANL Guest House, May 27-29, 2009

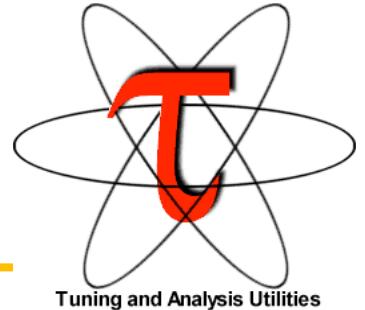
Sameer Shende and the TAU team

sameer@paratools.com

TAU: A Quick Reference



TAU Performance System

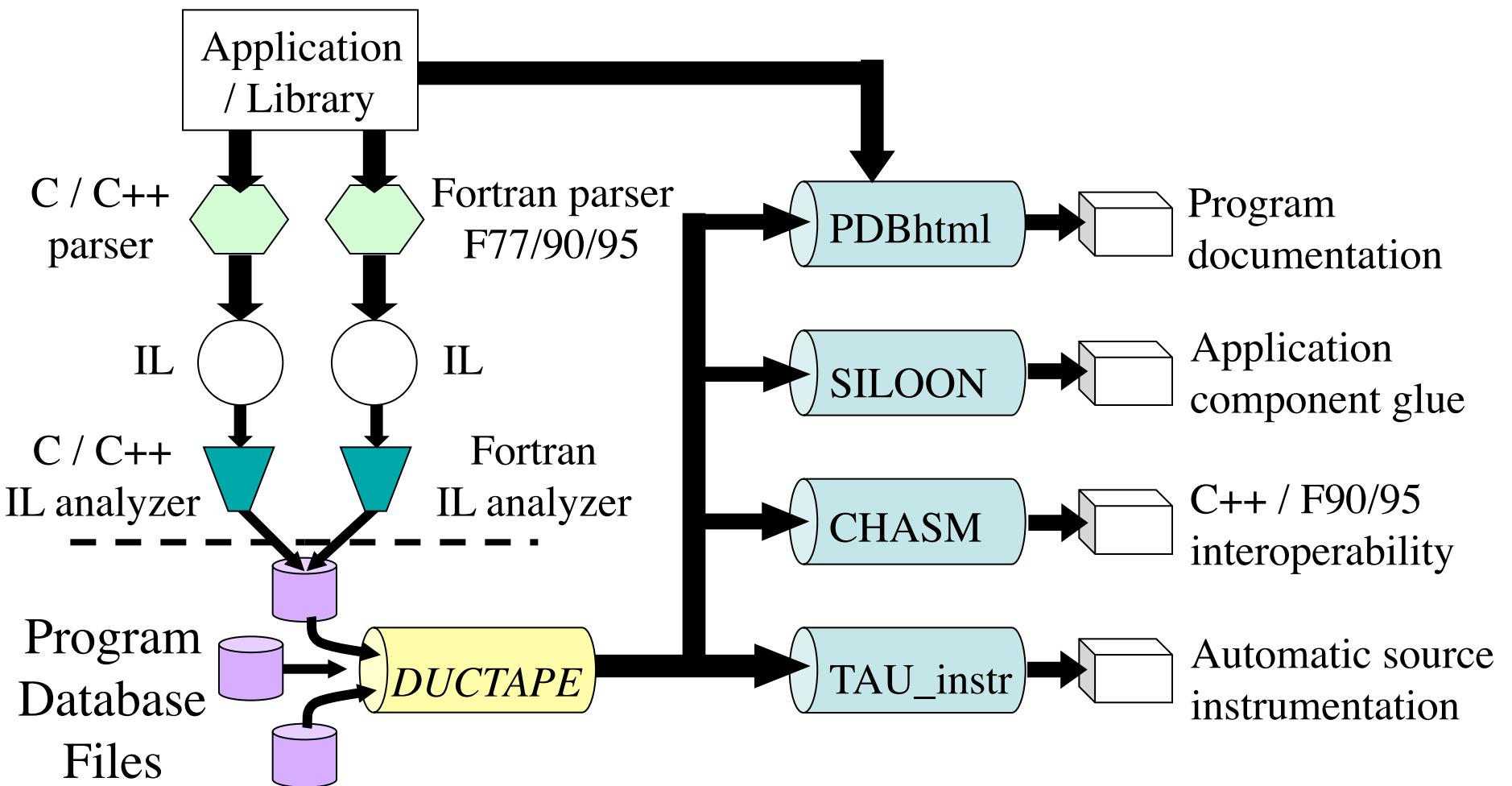


- <http://tau.uoregon.edu/>
- Multi-level performance instrumentation
 - Multi-language automatic source instrumentation
- Flexible and configurable performance measurement
- Widely-ported parallel performance profiling system
 - Computer system architectures and operating systems
 - Different programming languages and compilers
- Support for multiple parallel programming paradigms
 - Multi-threading, message passing, mixed-mode, hybrid
- Integration in complex software, systems, applications
- Partners: ANL, ORNL, LBL, LLNL, LANL, FZJ, UTK, TUD, PSC, UIUC

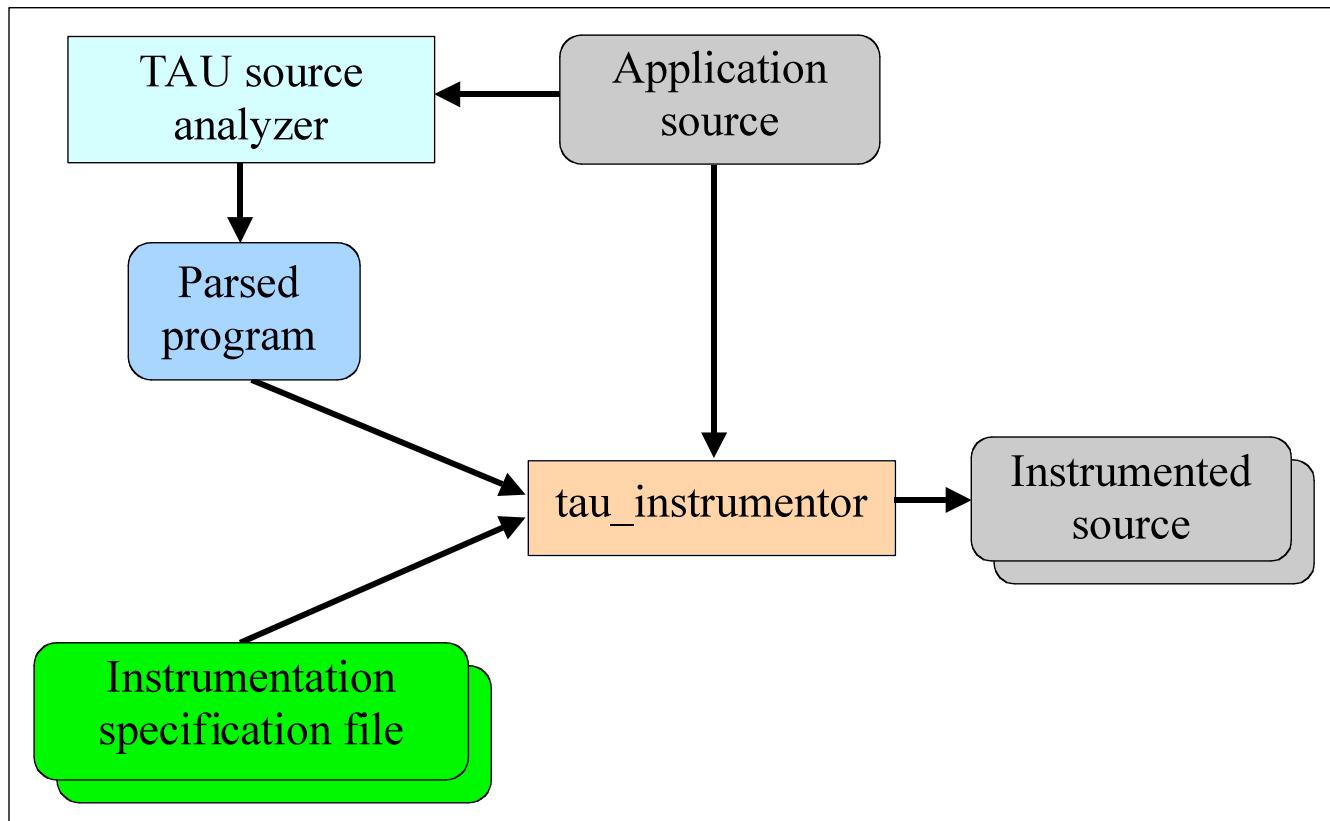
What is TAU?

- TAU is a performance evaluation tool
- It supports both parallel profiling and tracing
- Profiling shows you how much (total) time was spent in each routine (event)
- Tracing shows you *when* events take place in each process along a timeline
- Profiling and tracing can measure time as well as hardware performance counters
- TAU uses a package called *PDT* for automatic instrumentation of the source code
- With PDT, TAU can instrument routine, loop, phase, I/O, and memory
- TAU can also use your compiler to insert the instrumentation at routine boundaries
- TAU can *throttle* the insignificant lightweight routines at runtime to reduce perturbation. It can also *subtract* the timer overhead at runtime to compensate.
- TAU runs on all HPC platforms and it is free (BSD style license)
- TAU has instrumentation, measurement and analysis tools
 - ParaProf is TAU's 3D profile browser, PerfDMF is the TAU database tool
- To use TAU, all you need to do is set a couple of environment variables and substitute the name of your compiler with a TAU shell script (e.g., tau_f90.sh)

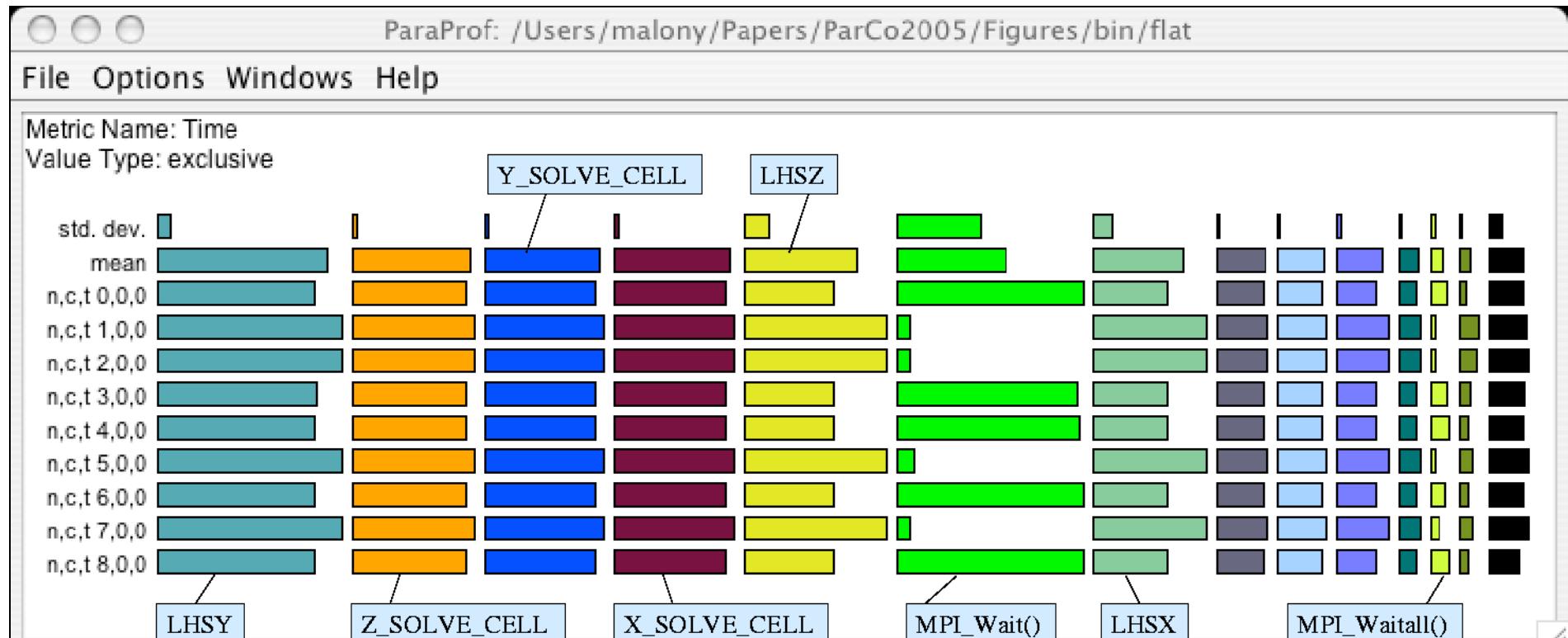
Program Database Toolkit (PDT)



Automatic Source-Level Instrumentation in TAU using Program Database Toolkit (PDT)



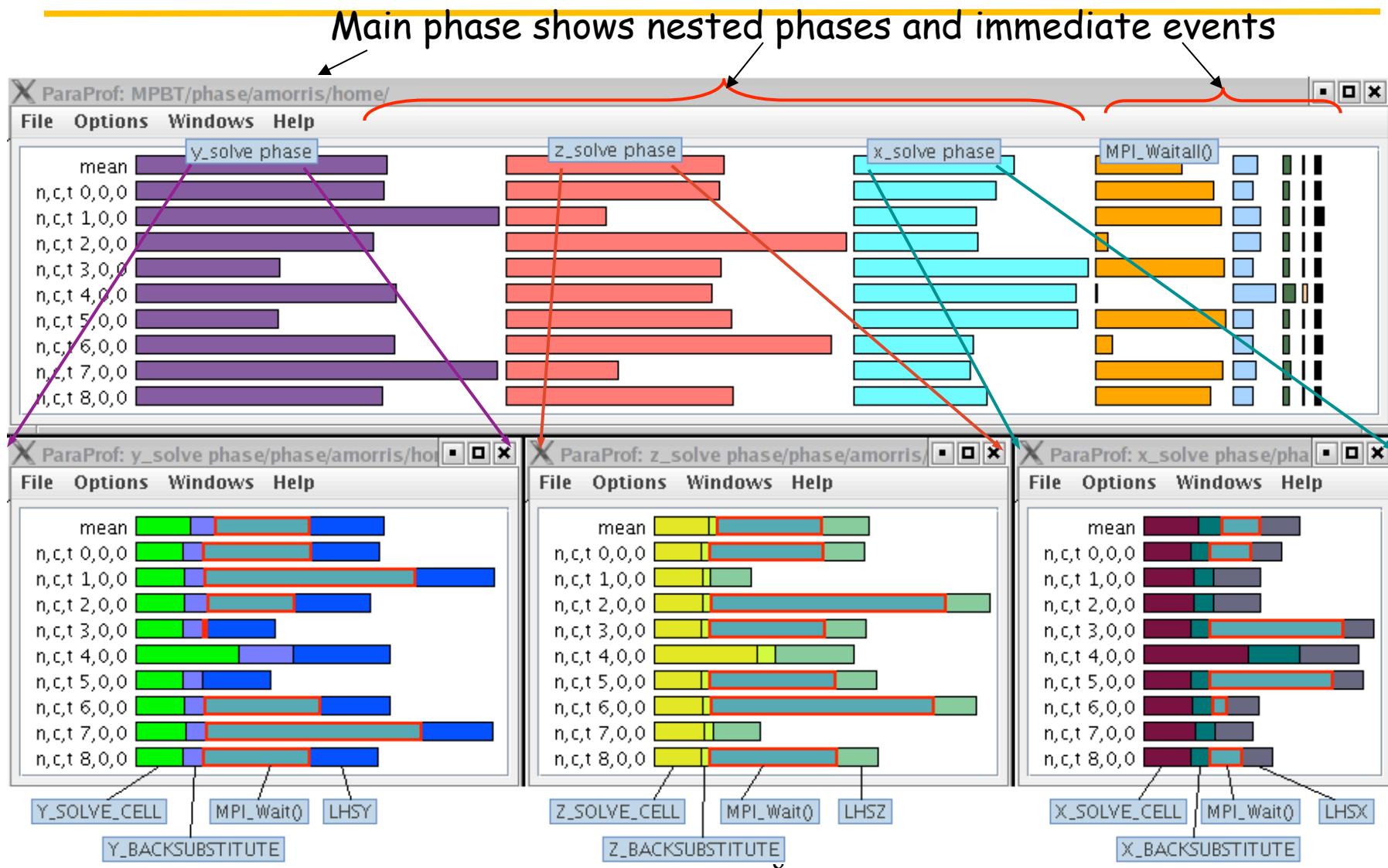
ParaProf – Flat Profile (NAS BT)



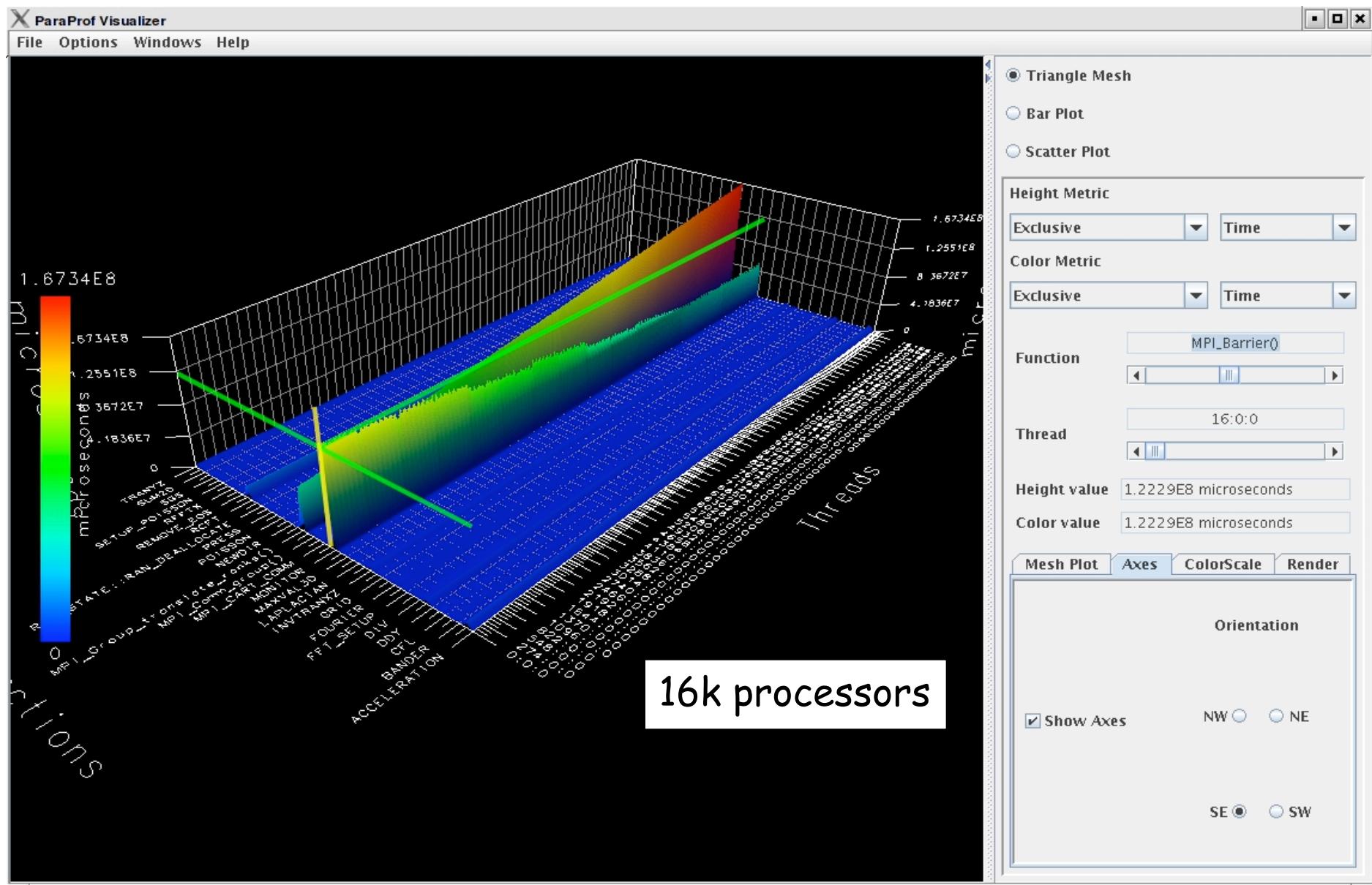
Application routine names
reflect phase semantics

How is MPI_Wait()
distributed relative to
solver direction?

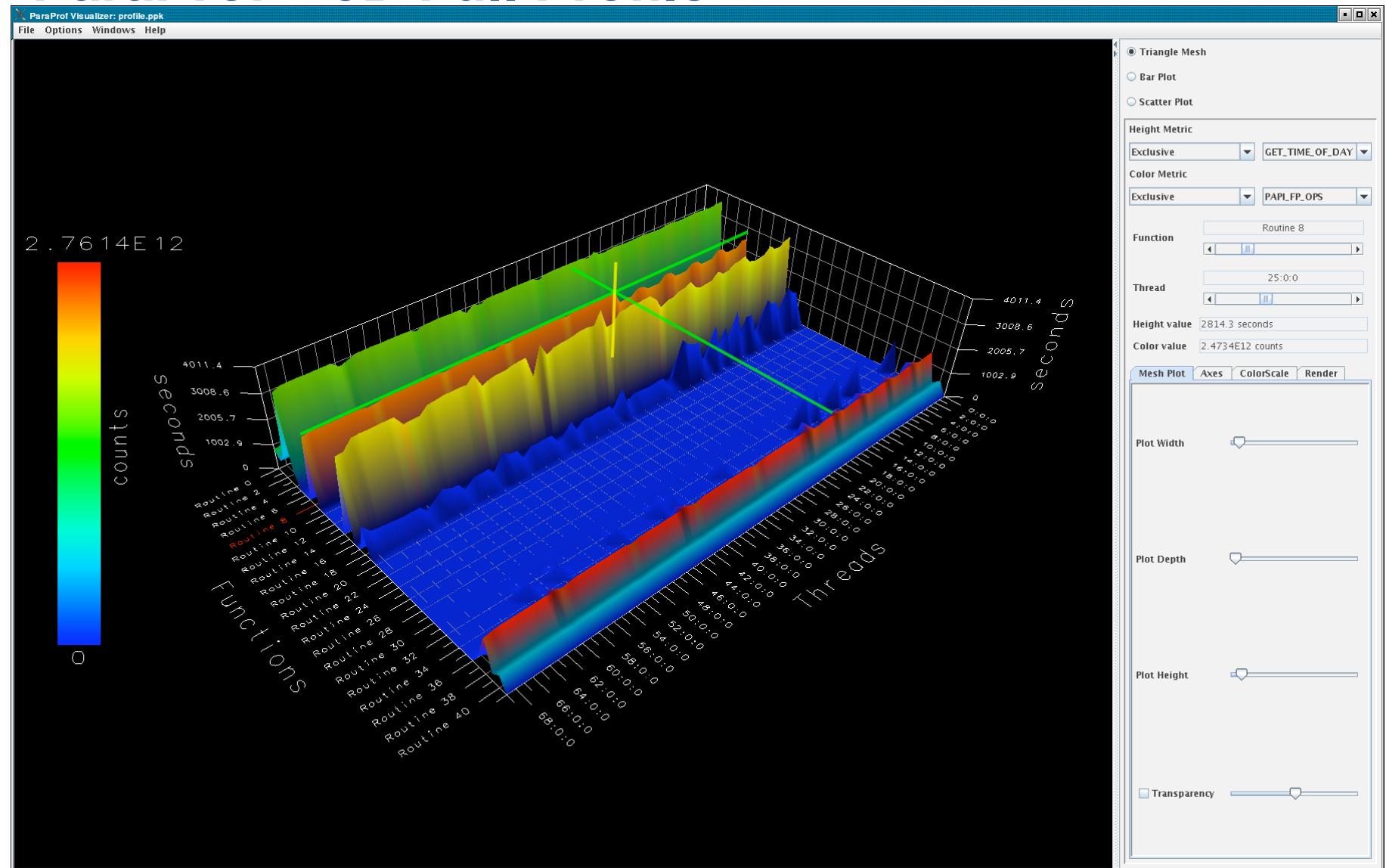
ParaProf – Phase Profile (NAS BT)



ParaProf – 3D Full Profile (Miranda, LLNL)



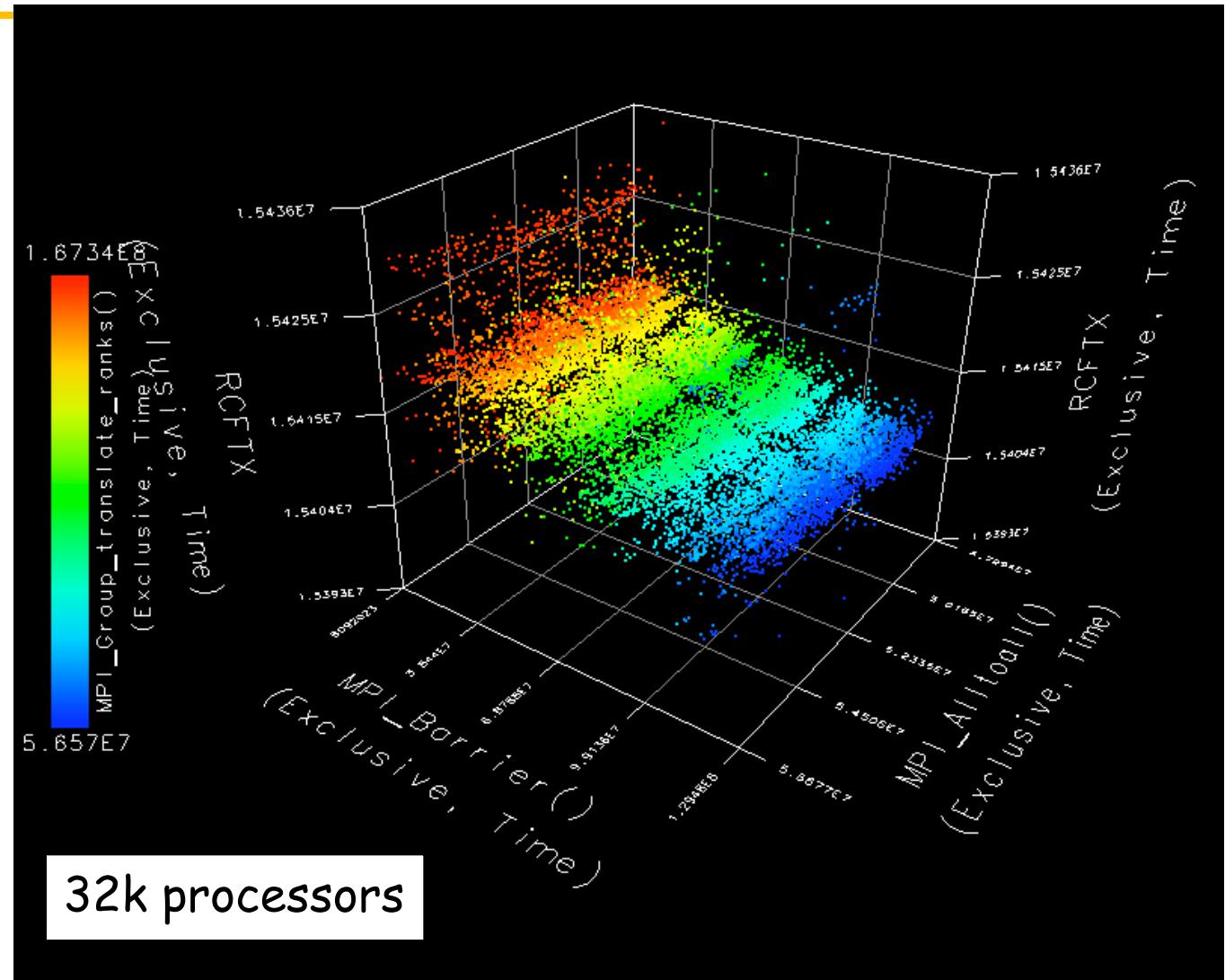
ParaProf – 3D Full Profile



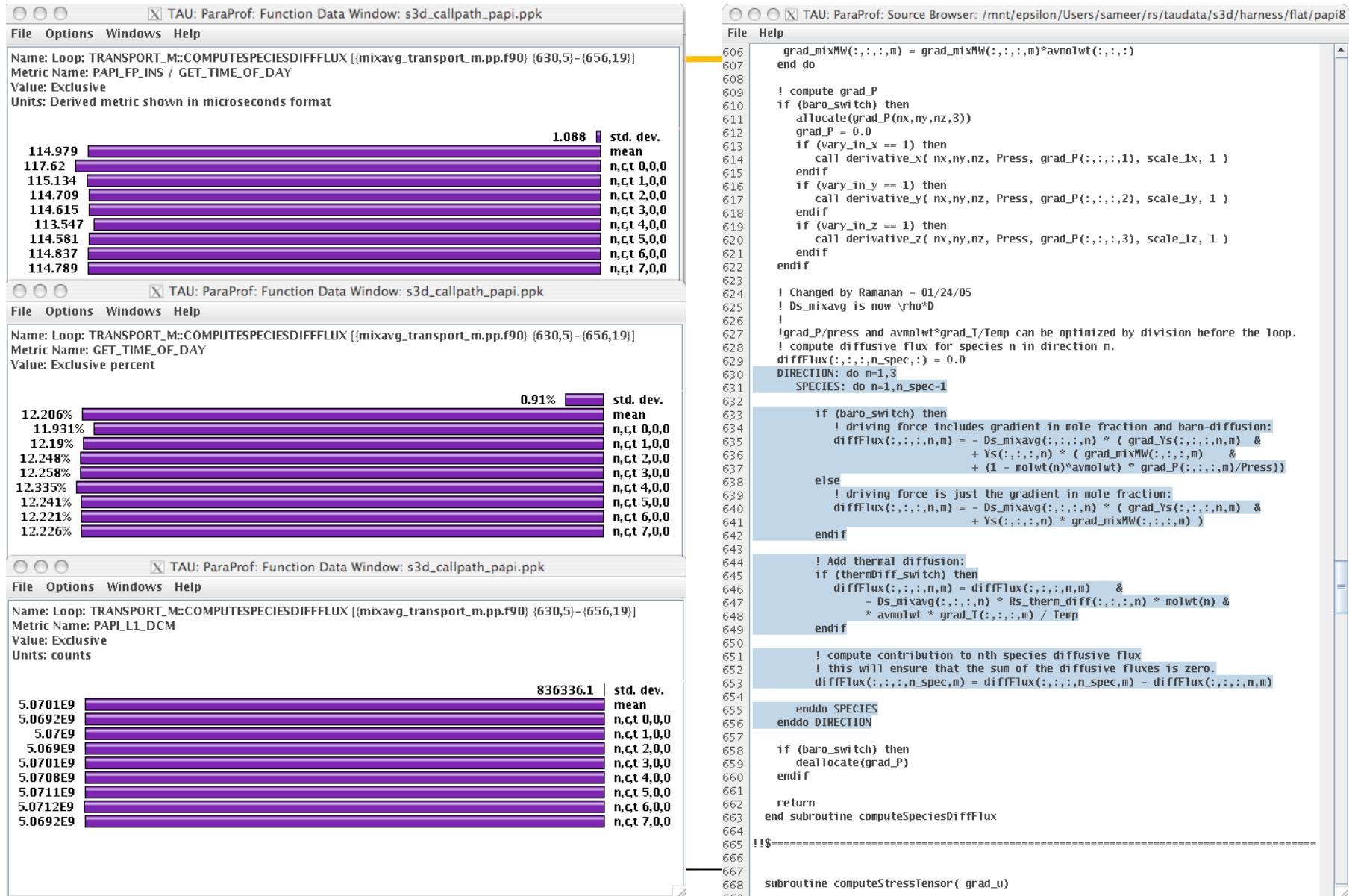
ParaTools

ParaProf – 3D Scatterplot

- Each point is a “thread” of execution
- A total of four metrics shown in relation
- 3D profile visualization library
 - JOGL



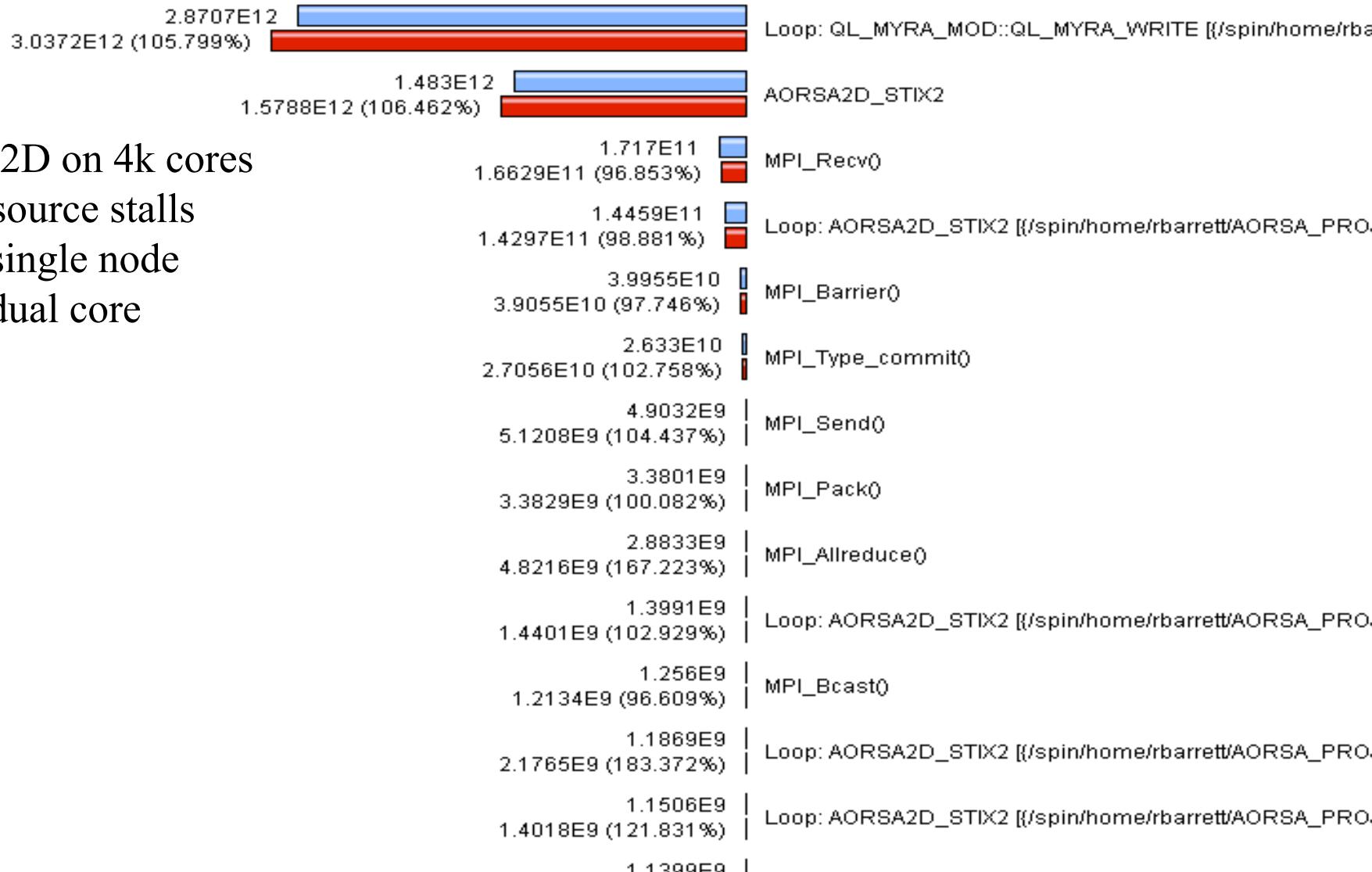
ParaProf's Source Browser: Loop Level Instrumentation



Comparing Effects of MultiCore Processors

Metric: PAPI_RES_STL
Value: Exclusive
Units: counts

C:\iter.350x350.4096pes.sn.loops.BARRIER.ppk - Mean
C:\iter.350x350.2048pes.dc.loops.BARRIER.ppk - Mean

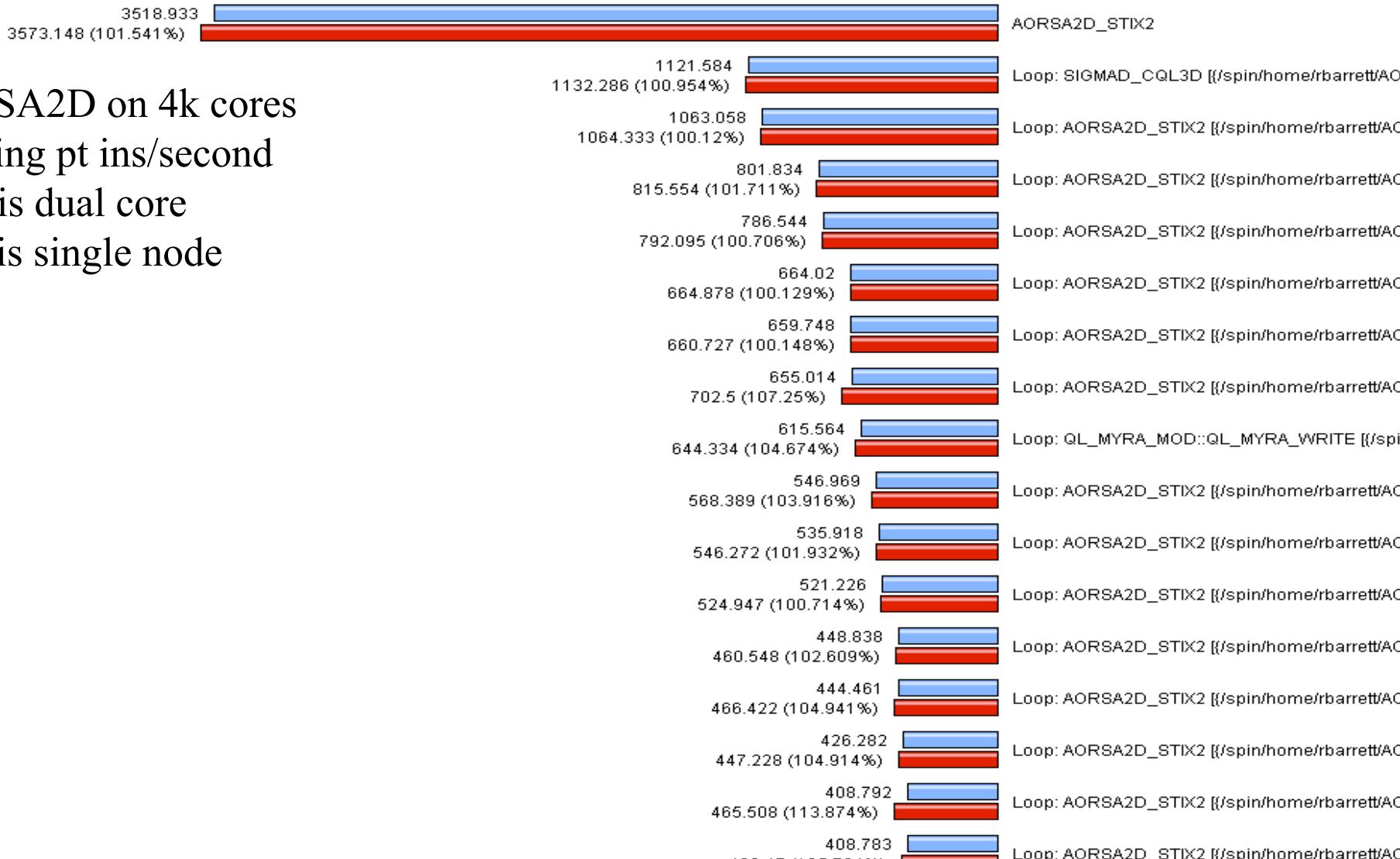


- AORSA2D on 4k cores
- PAPI resource stalls
- Blue is single node
- Red is dual core

Para

Comparing FLOPS: MultiCore Processors

Metric: PAPI_FP_OPS / GET_TIME_OF_DAY
Value: Exclusive
Units: Derived metric shown in microseconds format

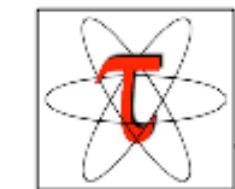


- AORSA2D on 4k cores
- Floating pt ins/second
- Blue is dual core
- Red is single node

Para

PerfDMF: Performance Data Mgmt. Framework

TAU Performance System



raw profiles

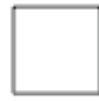


* gprof
* mpiP
* psrun
* HPMtoolkit
* ...

XML document



formatted profile data



profile metadata



Performance Analysis Programs

scalability analysis

ParaProf

cluster analysis

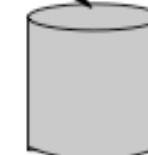


Query and Analysis Toolkit



Java PerfDMF API

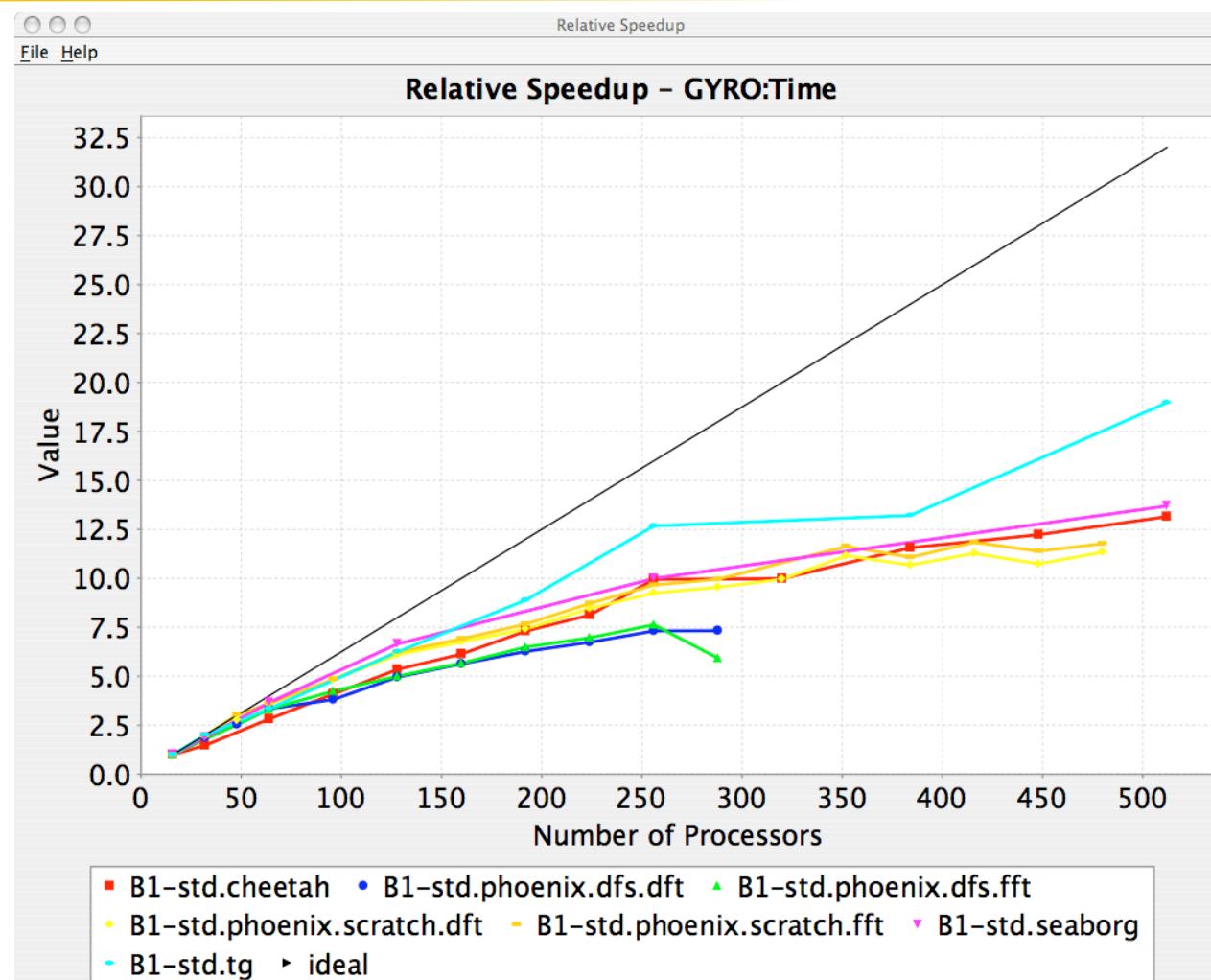
SQL (PostgreSQL, MySQL, DB2, Oracle)



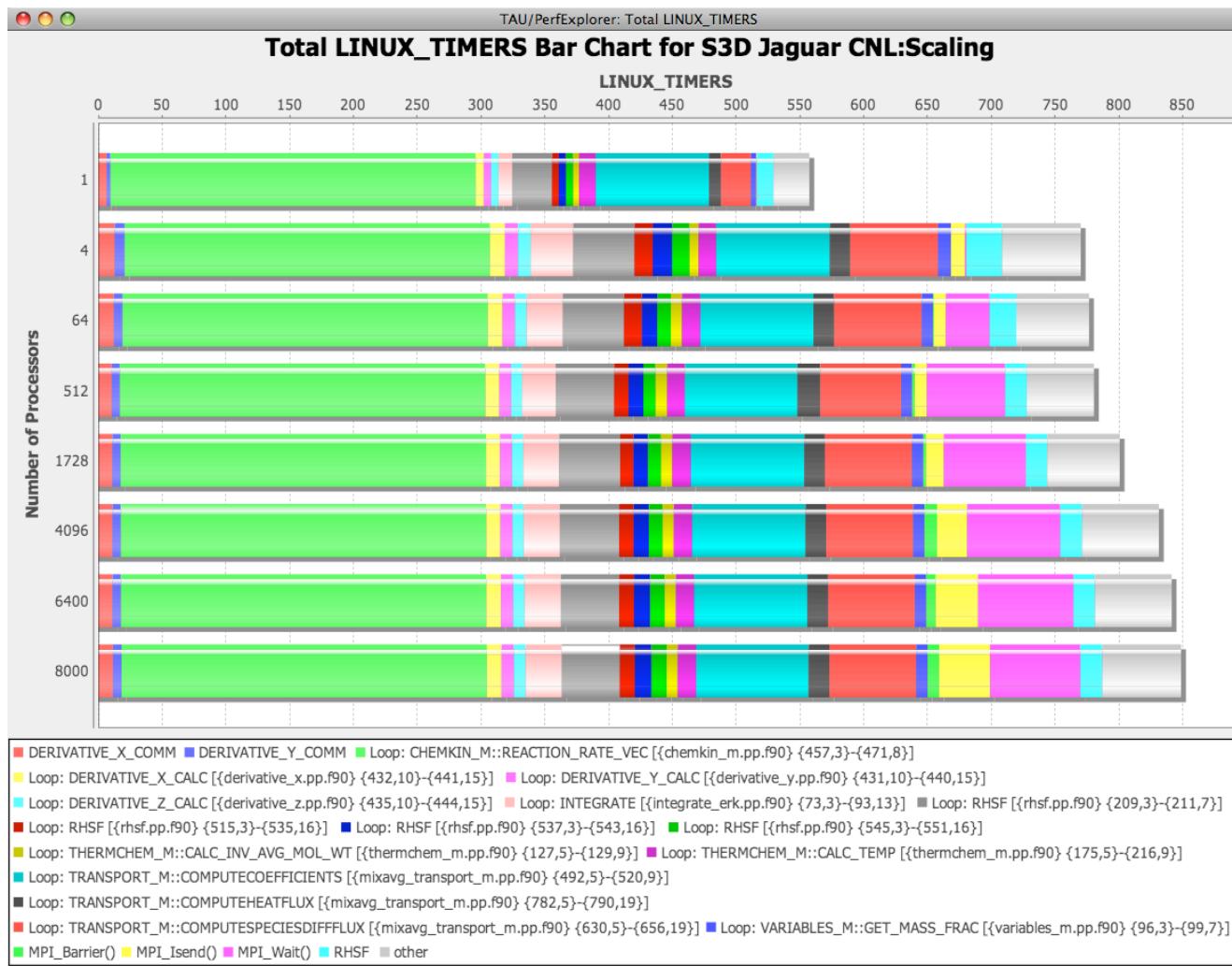
Data Mining
(Weka)

Statistics
(R / Omega)

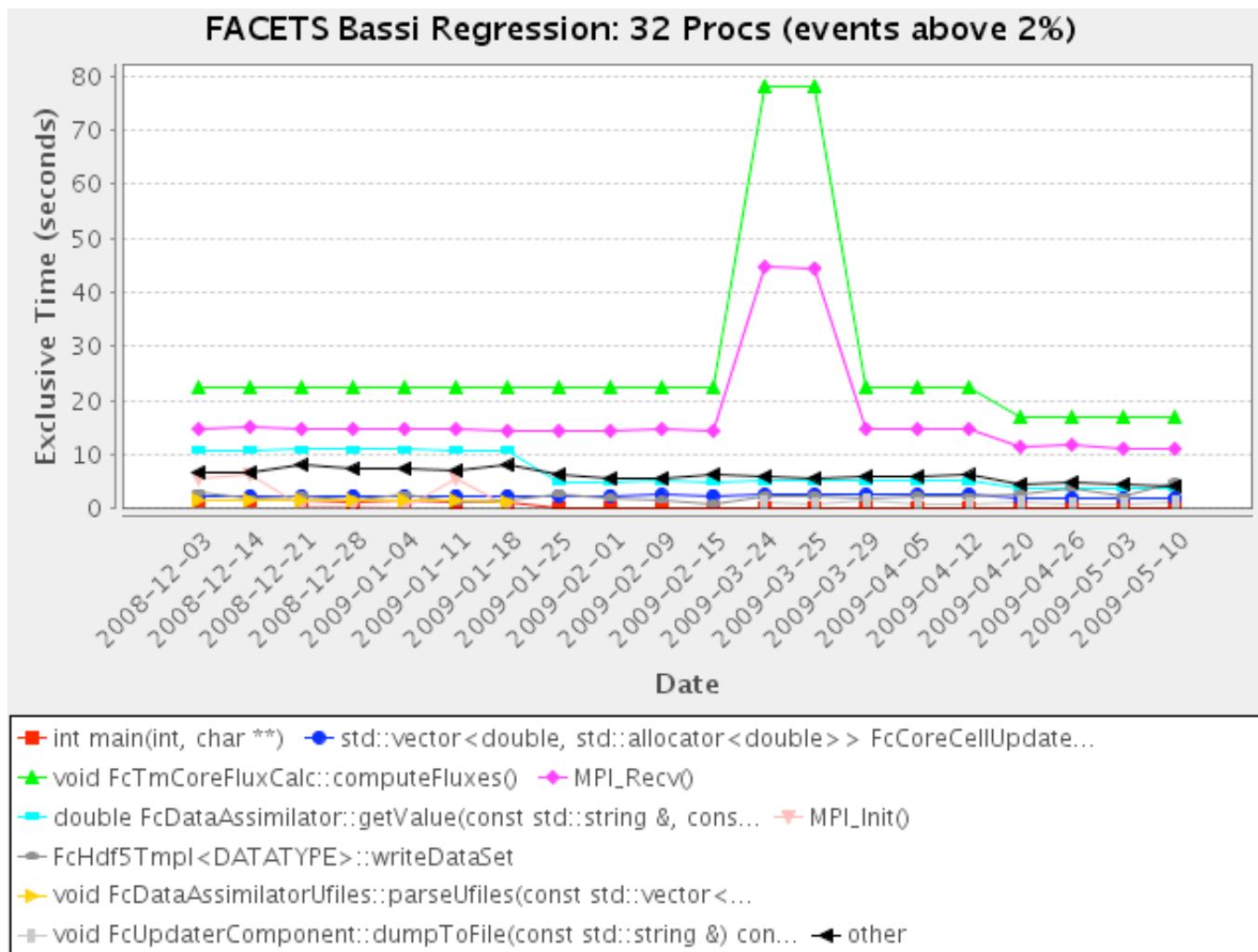
PerfExplorer: Comparing Relative Speedup on Different Architectures



Usage Scenarios: Evaluate Scalability



Performance Regression Testing



Profiling

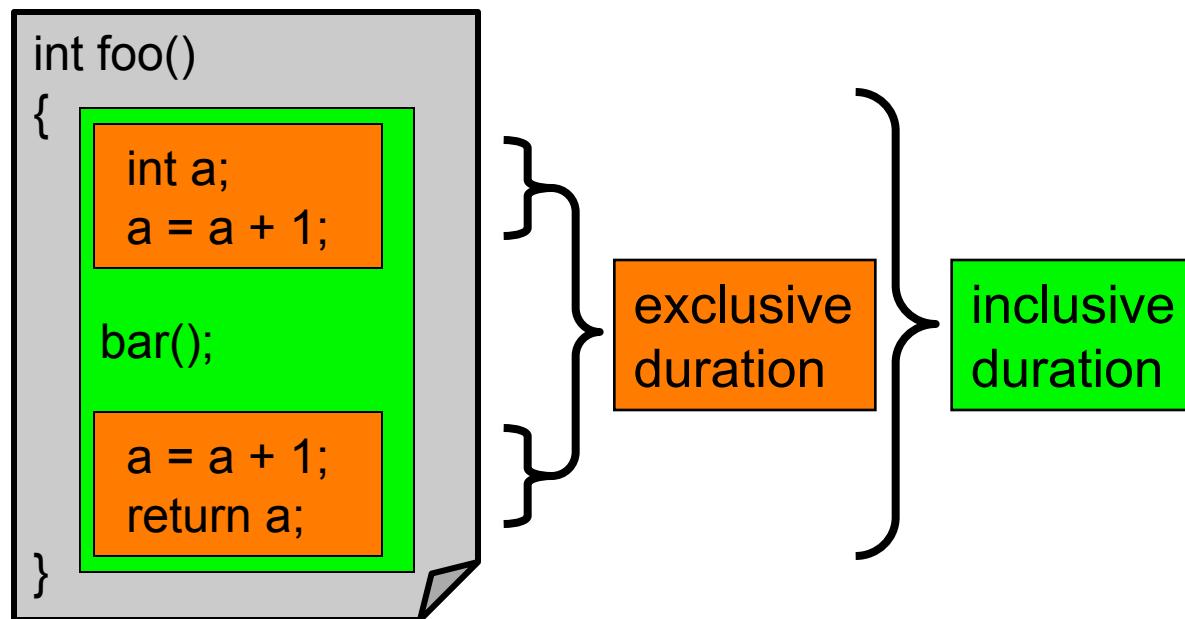
- Recording of aggregated information
 - Counts, time, ...
- ... about program and system entities
 - Functions, loops, basic blocks, ...
 - Processes, threads
- Methods
 - Event-based sampling (indirect, statistical)
 - OpenSpeedShop, PerfSuite, HPCToolkit, gprof,...
 - Direct measurement (deterministic)
 - TAU, VampirTrace, Scalasca,...

Direct Observation: Events

- Event types
 - Interval events (begin/end events)
 - measures performance between begin and end
 - metrics monotonically increase
 - Atomic events
 - used to capture performance data state
- Code events
 - Routines, classes, templates
 - Statement-level blocks, loops
- User-defined events
 - Specified by the user
- Abstract mapping events

Inclusive and Exclusive Profiles

- Performance with respect to code regions
- Exclusive measurements for region only
- Inclusive measurements includes child regions



Terminology – Example

- For routine “int main()”:
- Exclusive time
 - $100-20-50-20=10$ secs
- Inclusive time
 - 100 secs
- Calls
 - 1 call
- Subrs (no. of child routines called)
 - 3
- Inclusive time/call
 - 100secs

```
int main( )
{ /* takes 100 secs */

    f1(); /* takes 20 secs */
    f2(); /* takes 50 secs */
    f1(); /* takes 20 secs */

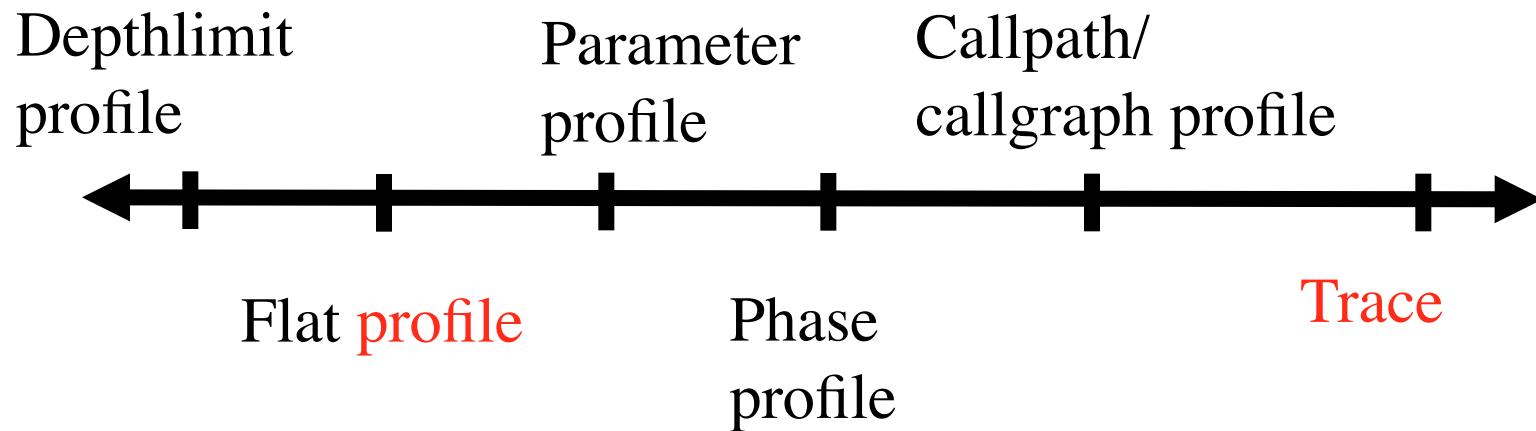
    /* other work */
}

/*
Time can be replaced by counts
from PAPI e.g., PAPI_FP_INS. */
```

Flat and Callpath Profiles

- Static call graph
 - Shows all parent-child calling relationships in a program
- Dynamic call graph
 - Reflects actual execution time calling relationships
- Flat profile
 - Performance metrics for when event is active
 - Exclusive and inclusive
- Callpath profile
 - Performance metrics for calling path (event chain)
 - Differentiate performance with respect to program execution state
 - Exclusive and inclusive

Performance Evaluation Alternatives



Each alternative has:

- one metric/counter
- multiple counters

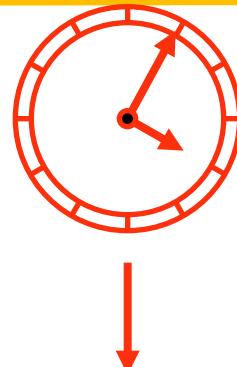
Volume of performance data



Tracing Measurement

Process A:

```
void master {  
    trace(ENTER, 1);  
    ...  
    trace(SEND, B);  
    send(B, tag, buf);  
    ...  
    trace(EXIT, 1);  
}
```



1	master
2	worker
3	...

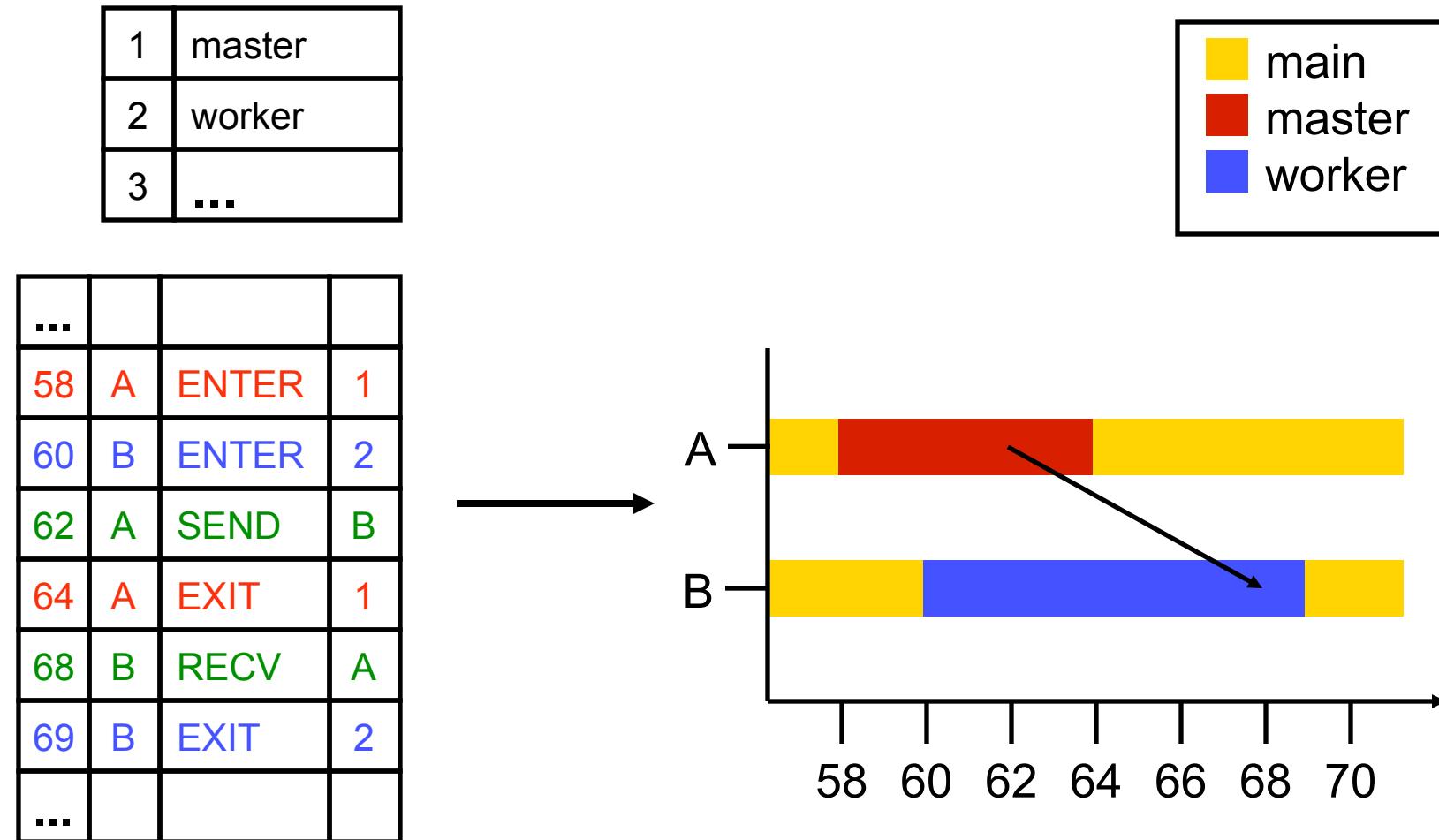
MONITOR

Process B:

```
void worker {  
    trace(ENTER, 2);  
    ...  
    recv(A, tag, buf);  
    trace(RECV, A);  
    ...  
    trace(EXIT, 2);  
}
```

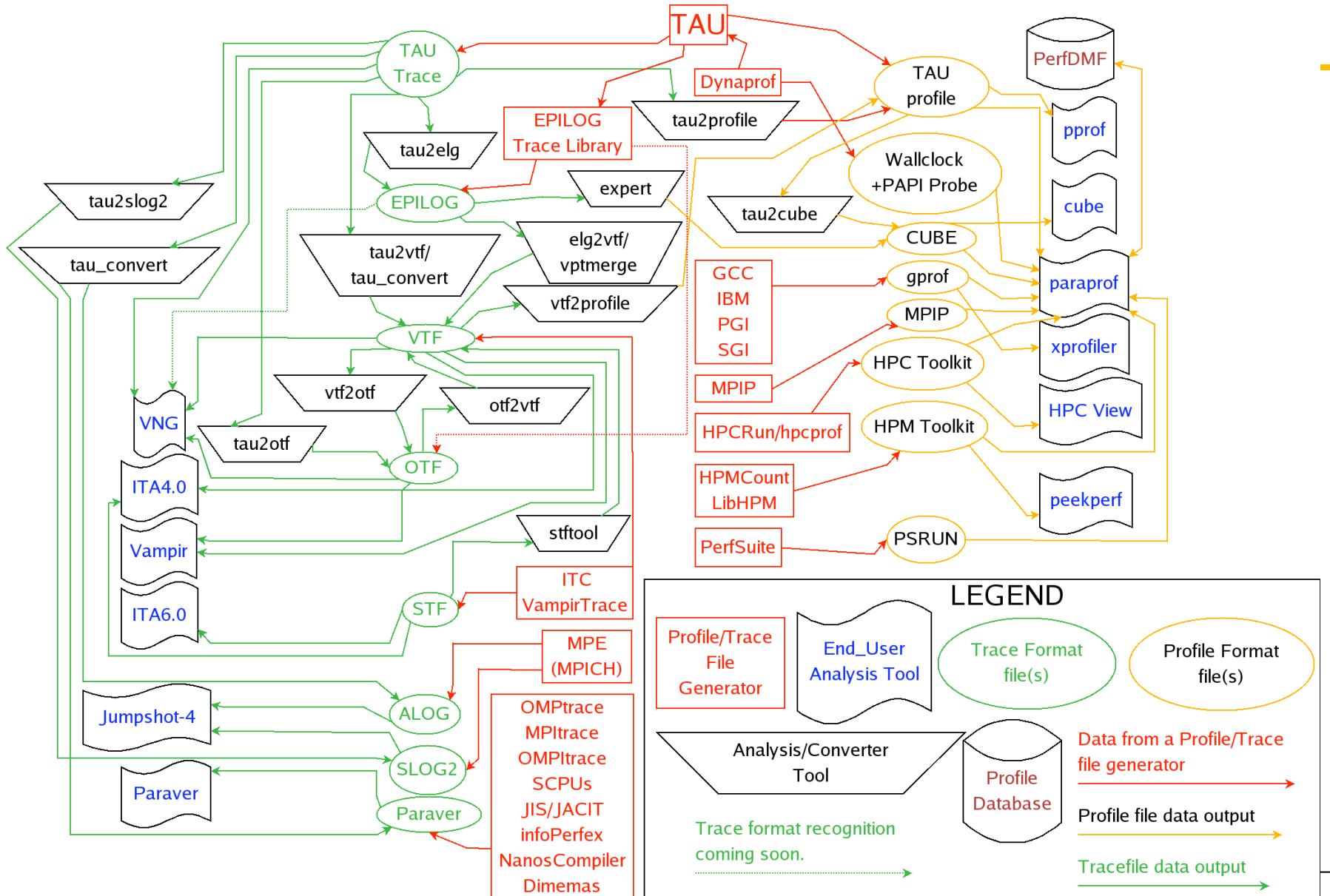
...				
58	A	ENTER	1	
60	B	ENTER	2	
62	A	SEND	B	
64	A	EXIT	1	
68	B	RECV	A	
69	B	EXIT	2	
...				

Tracing Analysis and Visualization





Building Bridges to Other Tools



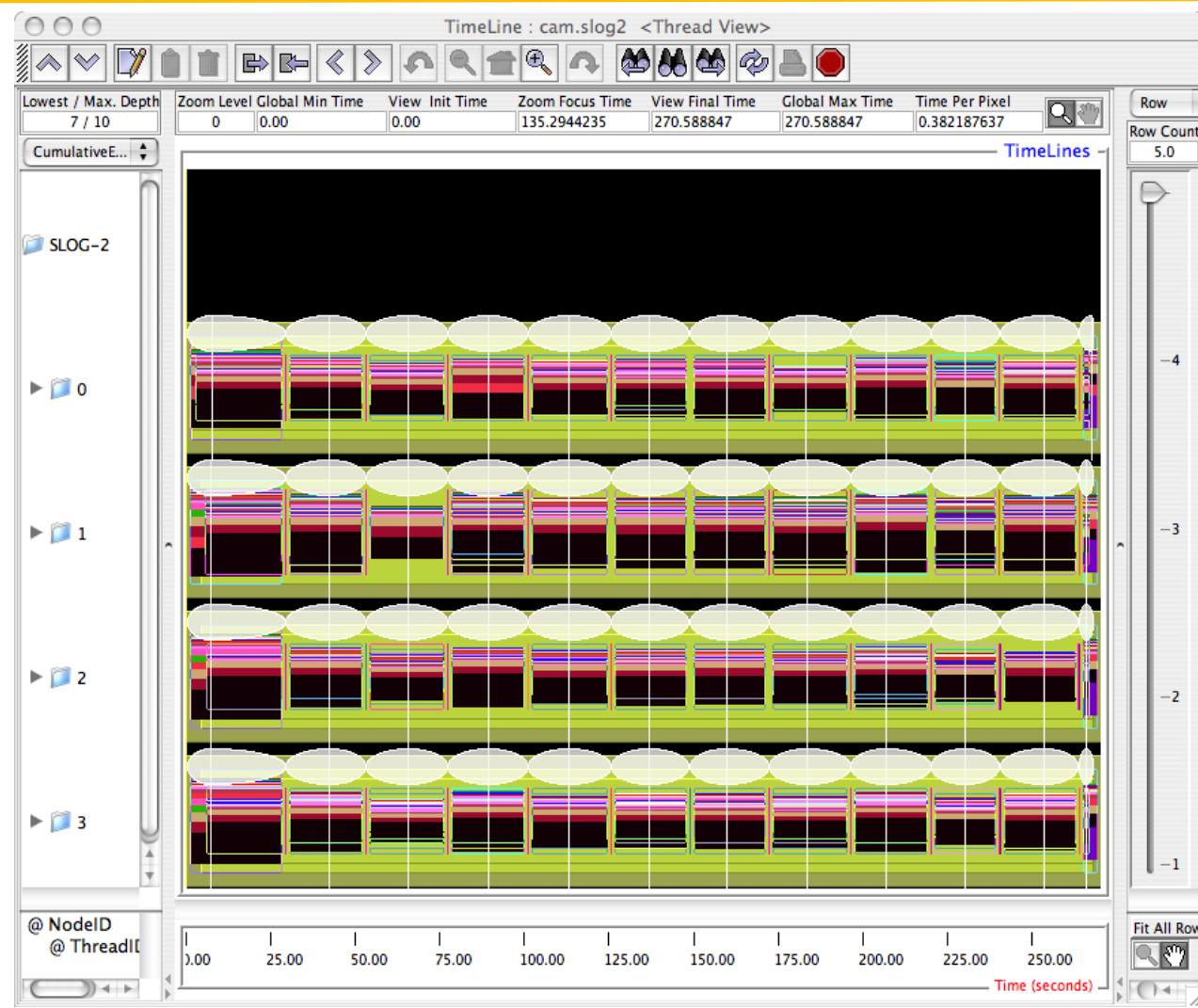
Trace Formats

- Different tools produce different formats
 - Differ by event types supported
 - Differ by ASCII and binary representations
 - Vampir Trace Format (VTF)
 - KOJAK (EPILOG)
 - Jumpshot (SLOG-2)
 - Paraver
- Open Trace Format (OTF)
 - Supports interoperation between tracing tools

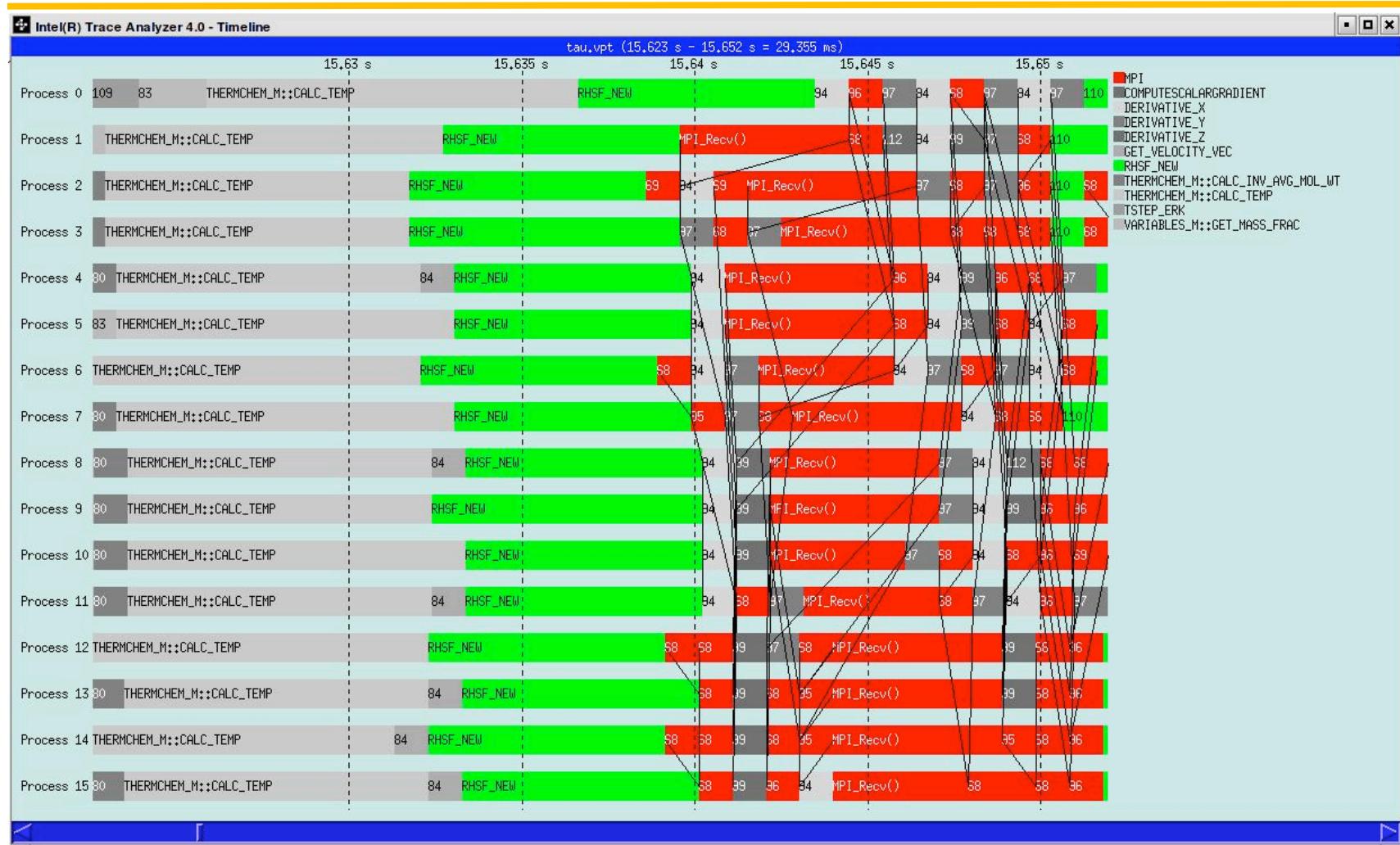
Profiling / Tracing Comparison

- Profiling
 - ☺ Finite, bounded performance data size
 - ☺ Applicable to both direct and indirect methods
 - ☹ Loses time dimension (not entirely)
 - ☹ Lacks ability to fully describe process interaction
- Tracing
 - ☺ Temporal and spatial dimension to performance data
 - ☺ Capture parallel dynamics and process interaction
 - ☹ Some inconsistencies with indirect methods
 - ☹ Unbounded performance data size (large)
 - ☹ Complex event buffering and clock synchronization

Jumpshot



Vampir – Trace Zoomed (S3D)



Steps of Performance Evaluation

- Collect basic routine-level timing profile to determine where most time is being spent
- Collect routine-level hardware counter data to determine types of performance problems
- Collect callpath profiles to determine sequence of events causing performance problems
- Conduct finer-grained profiling and/or tracing to pinpoint performance bottlenecks
 - Loop-level profiling with hardware counters
 - Tracing of communication operations

Using TAU: A brief Introduction

- TAU supports several measurement options (profiling, tracing, profiling with hardware counters, etc.)
- Each measurement configuration of TAU corresponds to a unique stub makefile that is generated when you configure it
- To instrument source code using PDT
 - Choose an appropriate TAU stub makefile in <arch>/lib:
`% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp/lib/Makefile.tau-mpi-pdt
% setenv TAU_OPTIONS '-optVerbose ...' (see tau_compiler.sh -help)`
And use tau_f90.sh, tau_cxx.sh or tau_cc.sh as Fortran, C++ or C compilers:
`% mpif90 foo.f90
changes to
% tau_f90.sh foo.f90`
- Execute application and analyze performance data:
`% pprof (for text based profile display)
% paraprof (for GUI)`

TAU Measurement Configuration

```
% cd /soft/apps/tau/tau-2.18.2/bgp/lib; ls Makefile.*  
Makefile.tau-pdt  
Makefile.tau-mpi-pdt  
Makefile.tau-opari-openmp-mpi-pdt  
Makefile.tau-mpi-scalasca-epilog-pdt  
Makefile.tau-mpi-vampirtrace-pdt  
Makefile.tau-mpi-papi-pdt  
Makefile.tau-papi-mpi-openmp-opari-pdt  
Makefile.tau-pthread-pdt...
```

- **For an MPI+F90 application, you may want to start with:**

Makefile.tau-mpi-pdt

- Supports MPI instrumentation & PDT for automatic source instrumentation
- % setenv TAU_MAKEFILE
/soft/apps/tau/tau-2.18.2/bgp/lib/Makefile.tau-mpi-pdt
- % tau_f90.sh matrix.f90 -o matrix

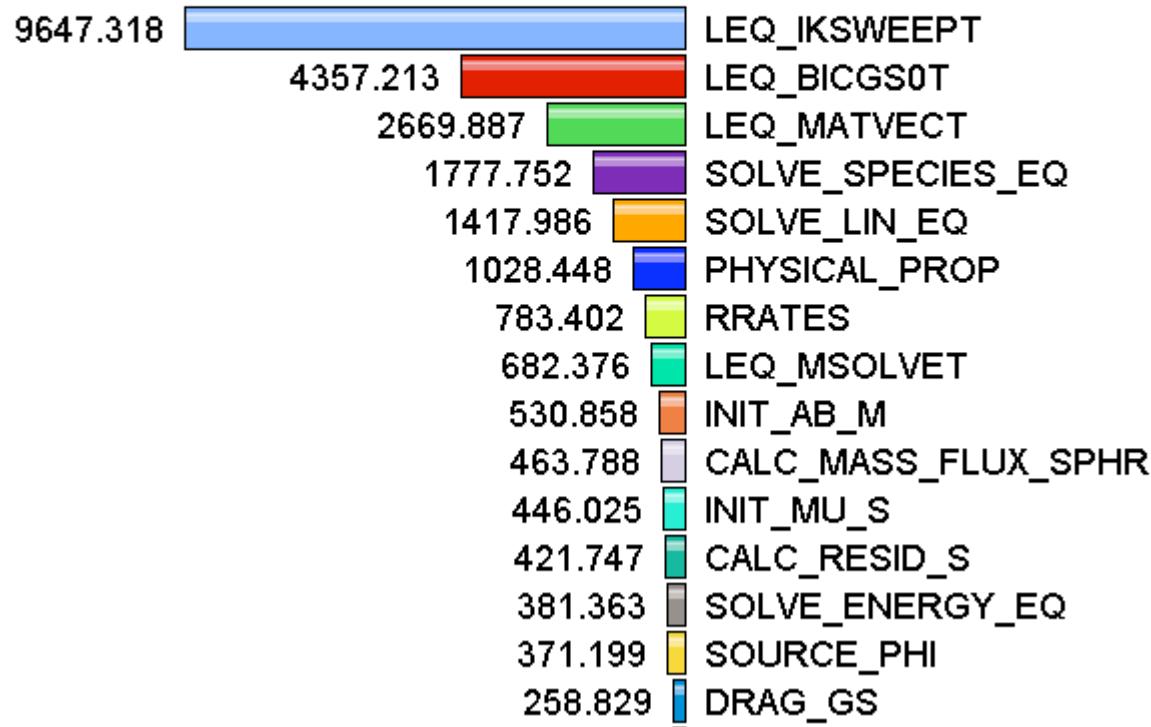
Usage Scenarios: Routine Level Profile

- Goal: What routines account for the most time? How much?
- Flat profile with wallclock time:

Metric: P_VIRTUAL_TIME

Value: Exclusive

Units: seconds



Solution: Generating a flat profile with MPI

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
                  /lib/Makefile.tau-mpi-pdt  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
OR  
% source /soft/apps/tau/src/tau.cshrc [ or tau.bashrc]  
% make F90=tau_f90.sh  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% qsub run.job  
% paraprof --pack app.ppk  
Move the app.ppk file to your desktop.  
  
% paraprof app.ppk
```

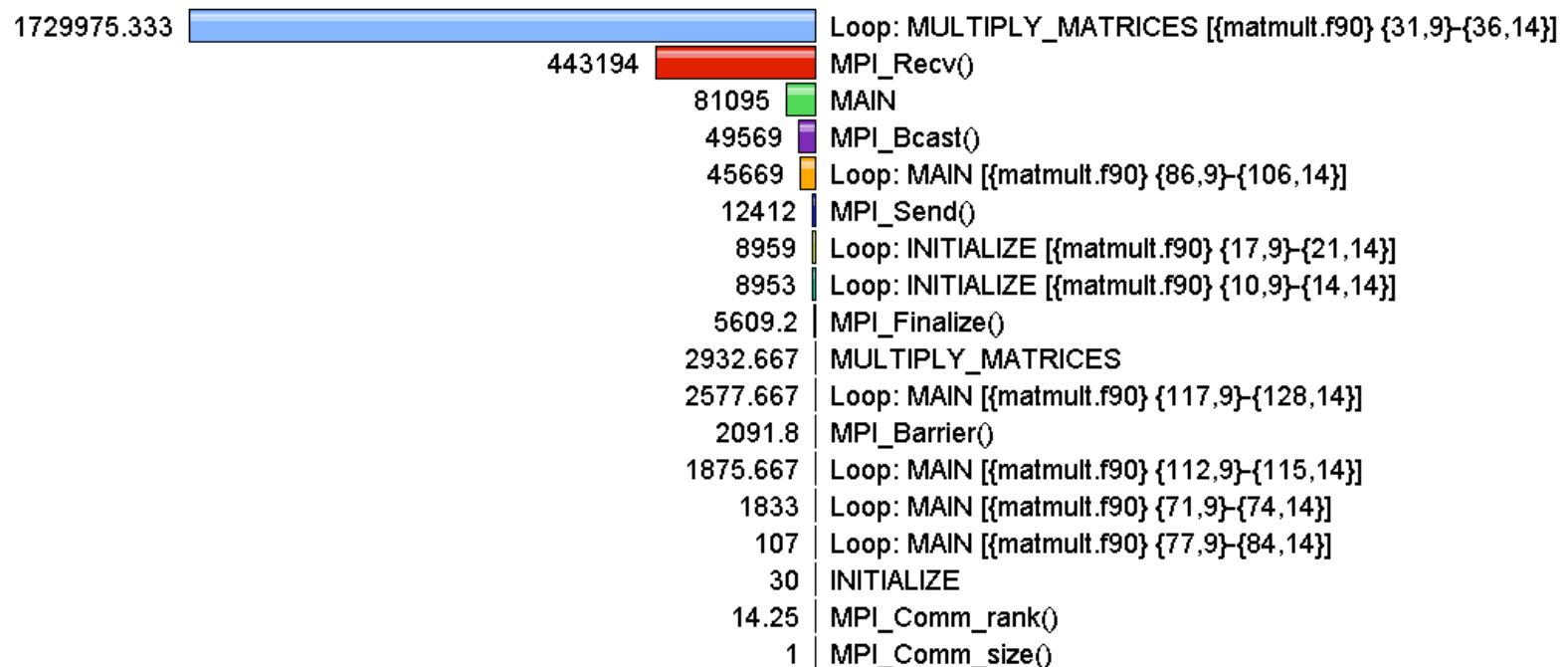
Usage Scenarios: Loop Level Instrumentation

- Goal: What loops account for the most time? How much?
- Flat profile with wallclock time with loop instrumentation:

Metric: GET_TIME_OF_DAY

Value: Exclusive

Units: microseconds



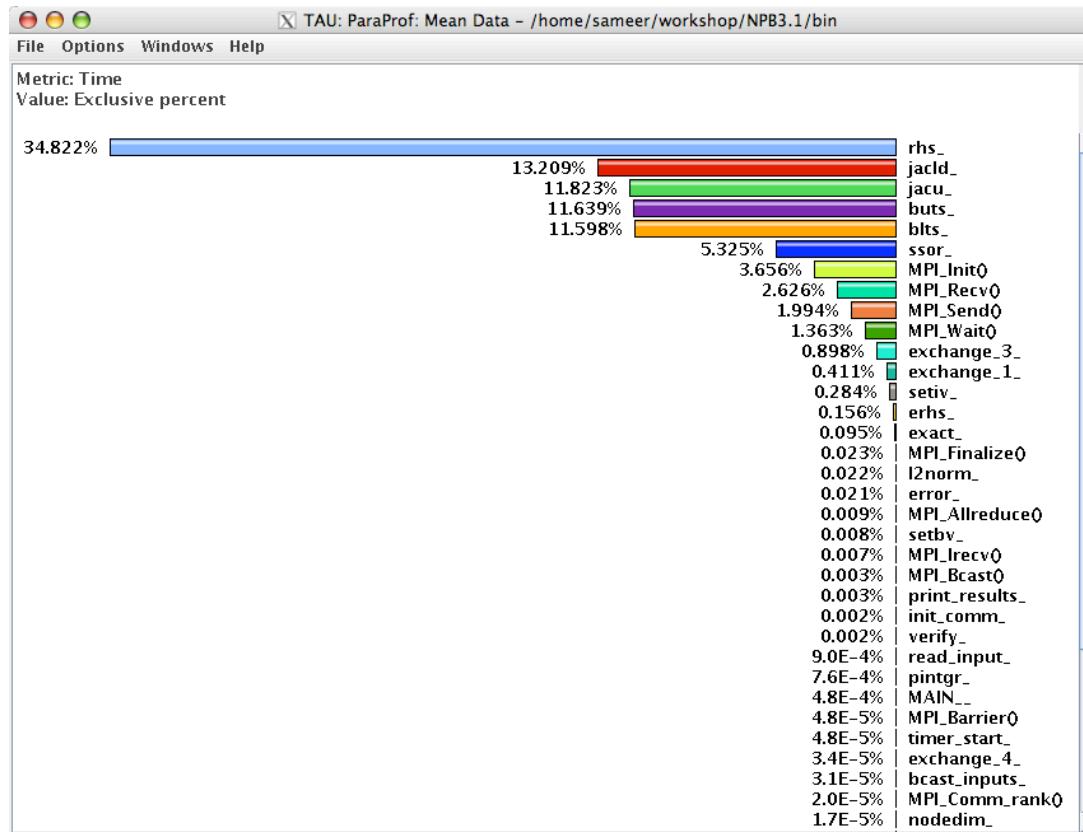
Par

Solution: Generating a loop level profile

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
                  /lib/Makefile.tau-mpi-pdt  
% setenv TAU_OPTIONS '-optTauSelectFile=select.tau -optVerbose'  
% cat select.tau  
BEGIN_INSTRUMENT_SECTION  
loops routine="#"  
END_INSTRUMENT_SECTION  
  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
% make F90=tau_f90.sh  
(Or edit Makefile and change F90=tau_f90.sh)  
% qsub run.job  
% paraprof --pack app.ppk  
Move the app.ppk file to your desktop.  
  
% paraprof app.ppk
```

Usage Scenarios: Compiler-based Instrumentation

- Goal: Easily generate routine level performance data using the compiler instead of PDT for parsing the source code



Use Compiler-Based Instrumentation

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
      /lib/Makefile.tau-mpi  
  
% setenv TAU_OPTIONS '-optCompInst -optVerbose'  
  
% % set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
% make F90=tau_f90.sh  
  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% qsub run.job  
% paraprof --pack app.ppk  
  Move the app.ppk file to your desktop.  
% paraprof app.ppk
```

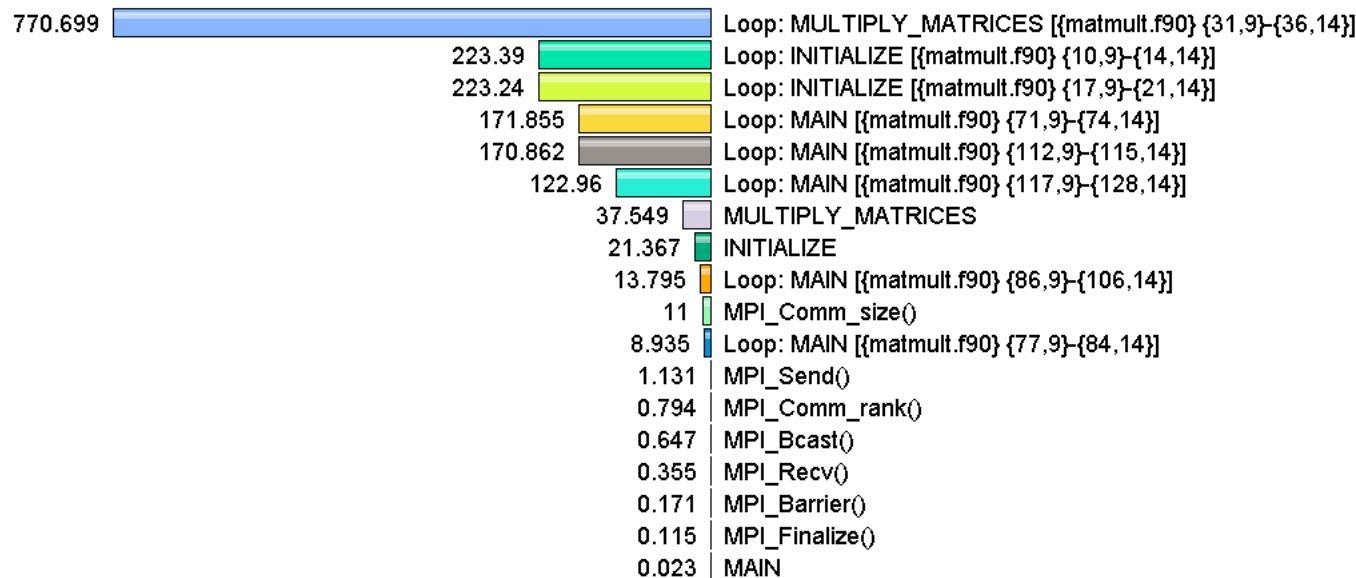
Usage Scenarios: Calculate mflops in Loops

- Goal: What MFlops am I getting in all loops?
- Flat profile with PAPI_FP_INS/OPS and time (-multiplecounters) with loop instrumentation:

Metric: PAPI_FP_INS / GET_TIME_OF_DAY

Value: Exclusive

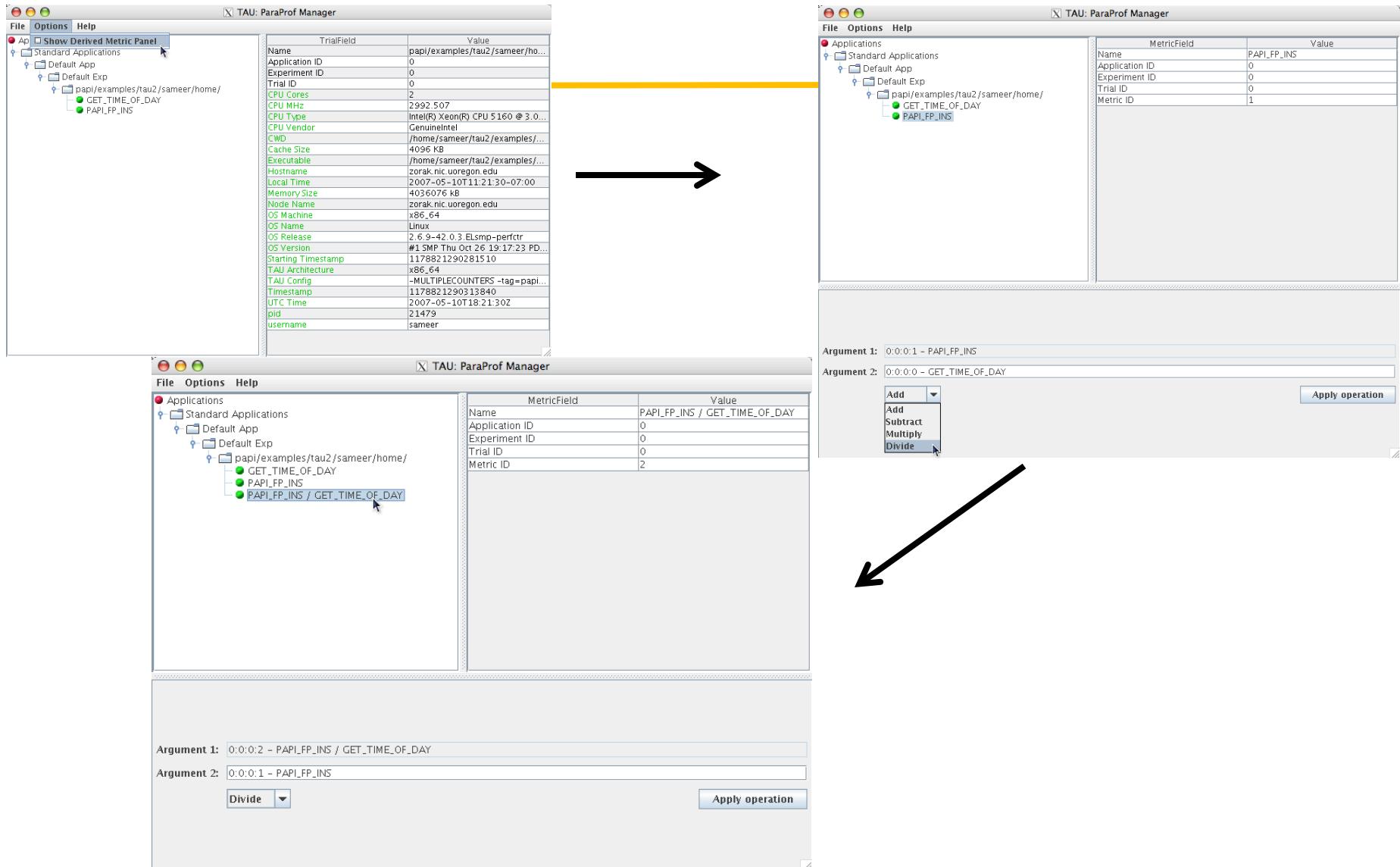
Units: Derived metric shown in microseconds format



Generate a PAPI profile with 2 or more counters

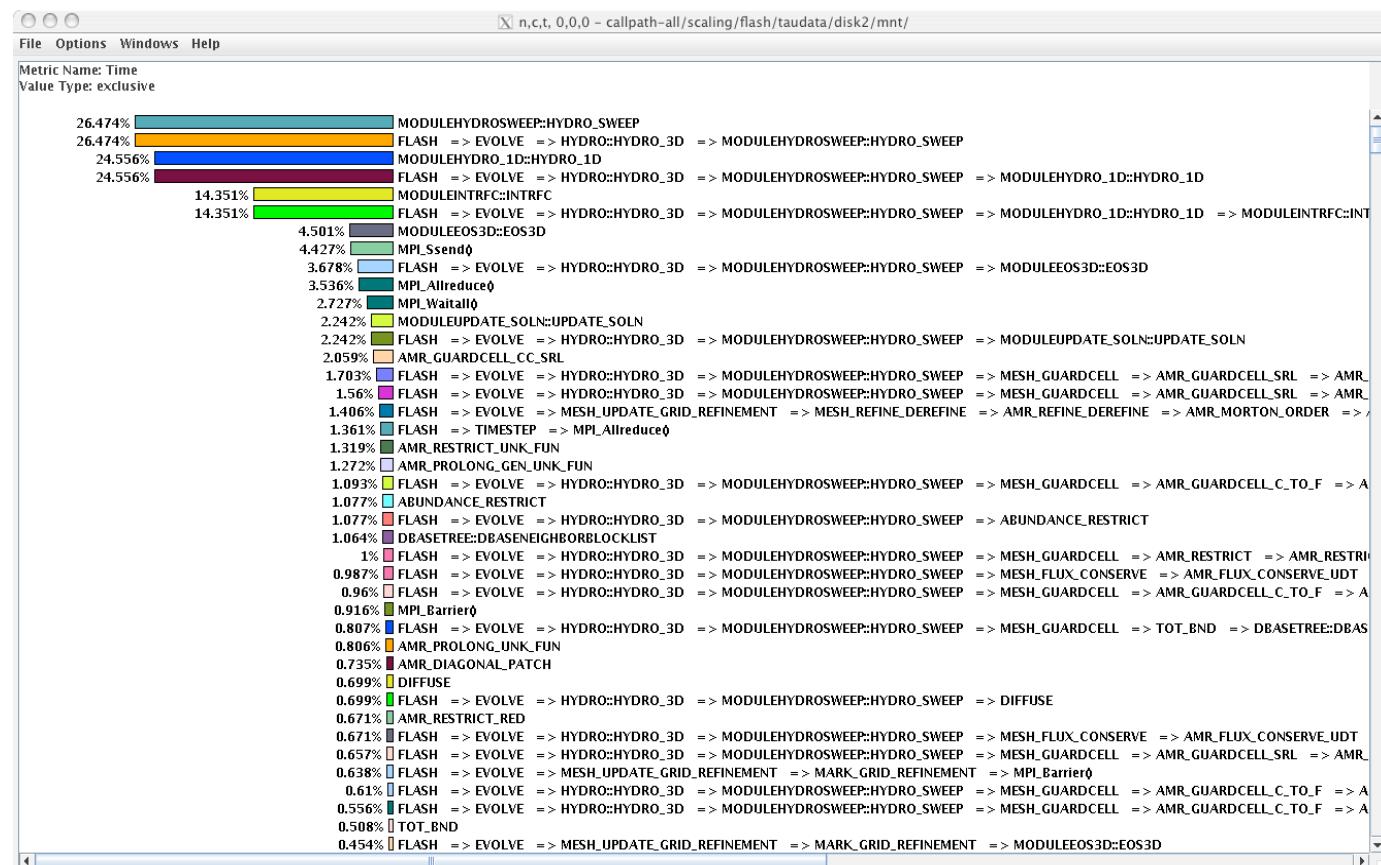
```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
          /lib/Makefile.tau-papi-mpi-pdt  
  
% setenv TAU_OPTIONS '-optTauSelectFile=select.tau -optVerbose'  
  
% cat select.tau  
  
BEGIN_INSTRUMENT_SECTION  
loops routine="#"  
END_INSTRUMENT_SECTION  
  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
% make F90=tau_f90.sh  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% setenv COUNTER1 GET_TIME_OF_DAY  
% setenv COUNTER2 PAPI_FP_INS  
  
OR  
  
% setenv TAU_METRICS TIME:PAPI_FP_INS  
% setenv TAU_METRICS TIME:PAPI_NATIVE_<event_name>:....  
% qsub run.job  
% paraprof --pack app.ppk           Move the app.ppk file to your desktop.  
% paraprof app.ppk  
  
Choose Options -> Show Derived Panel -> Arg 1 = PAPI_FP_INS,  
Arg 2 = GET_TIME_OF_DAY, Operation = Divide -> Apply, choose.
```

Derived Metrics in ParaProf



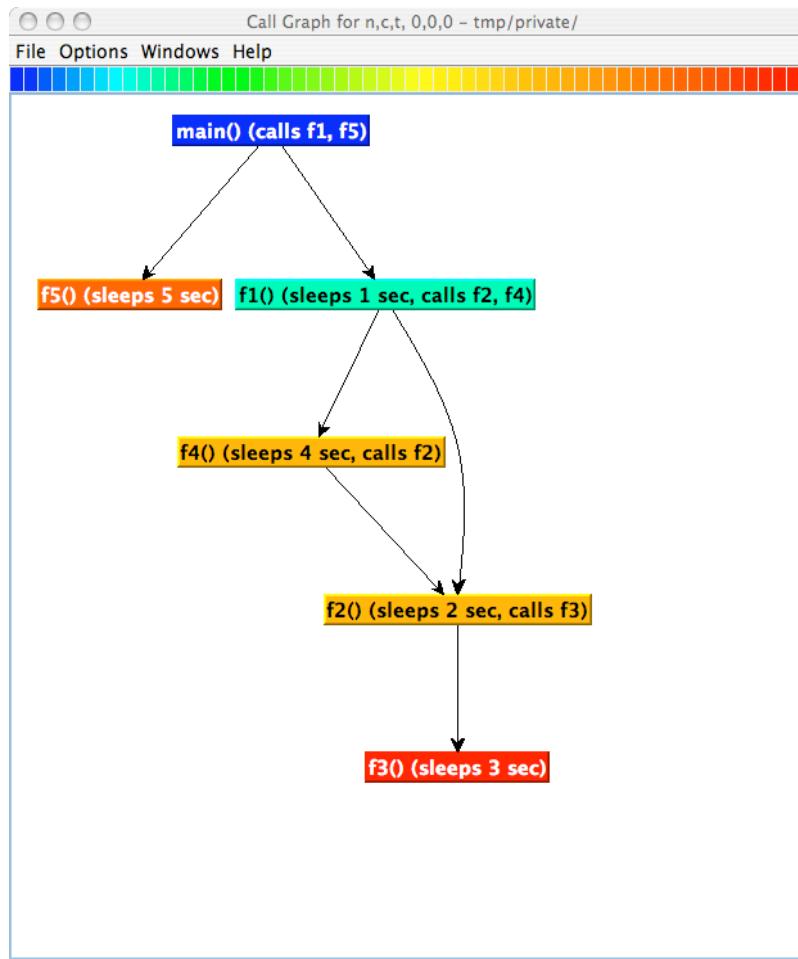
Usage Scenarios: Generating Callpath Profile

- Goal: Who calls my MPI_Barrier()? Where?
- Callpath profile for a given callpath depth:



Callpath Profile

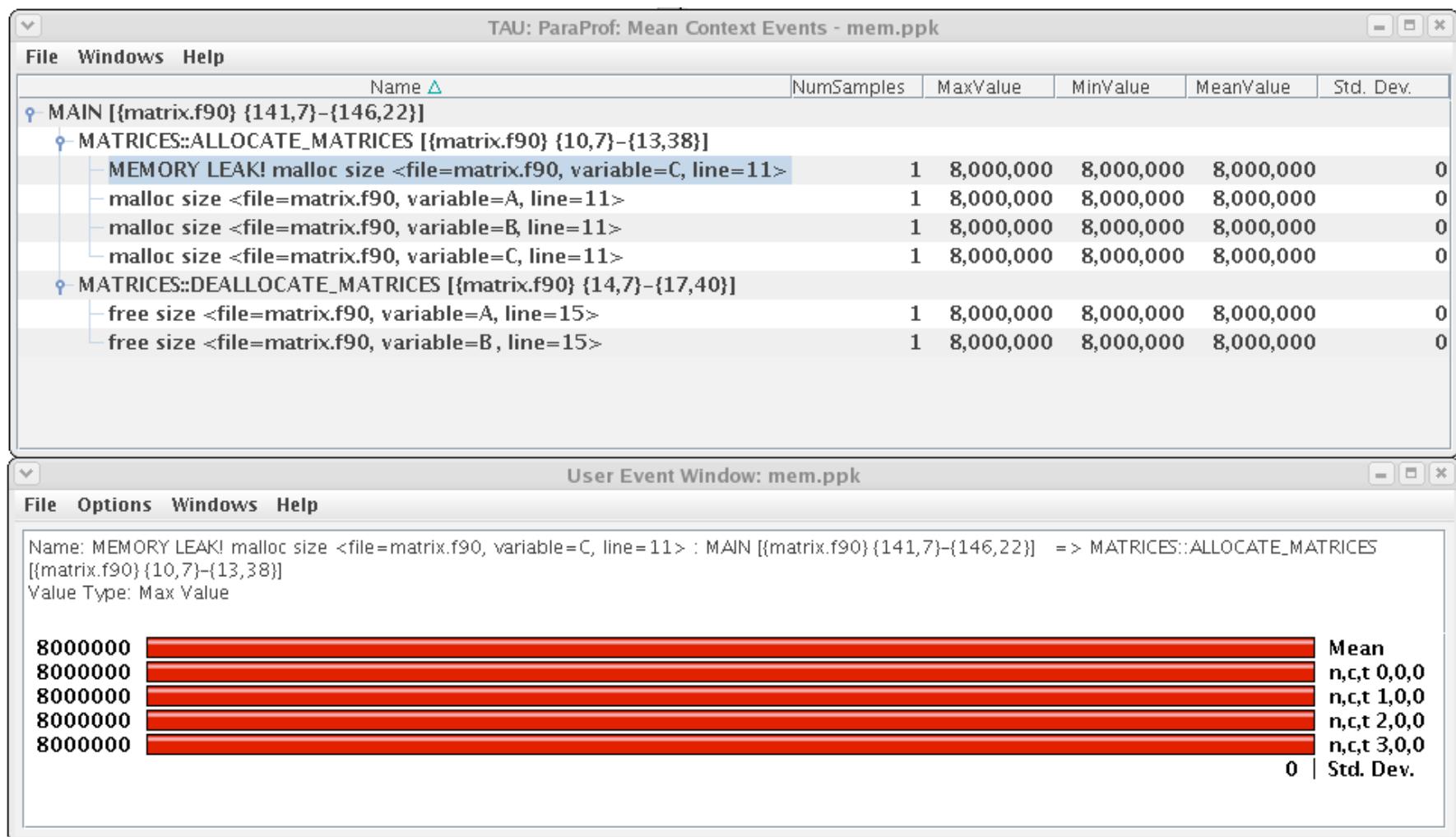
- Generates program callgraph



Generate a Callpath Profile

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
                  /lib/Makefile.tau-callpath-mpi-pdt  
  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
  
% make F90=tau_f90.sh  
  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% setenv TAU_CALLPATH_DEPTH 100  
  
  
% qsub run.job  
  
% paraprof --pack app.ppk  
  
  Move the app.ppk file to your desktop.  
  
% paraprof app.ppk  
  
(Windows -> Thread -> Call Graph)  
  
  
NOTE: In TAU v2.18.1+, you may choose to just set:  
  
% setenv TAU_CALLPATH 1  
  
instead of recompiling your code with the above stub makefile.  
Any TAU instrumented executable can generate callpath profiles.
```

Usage Scenario: Detect Memory Leaks

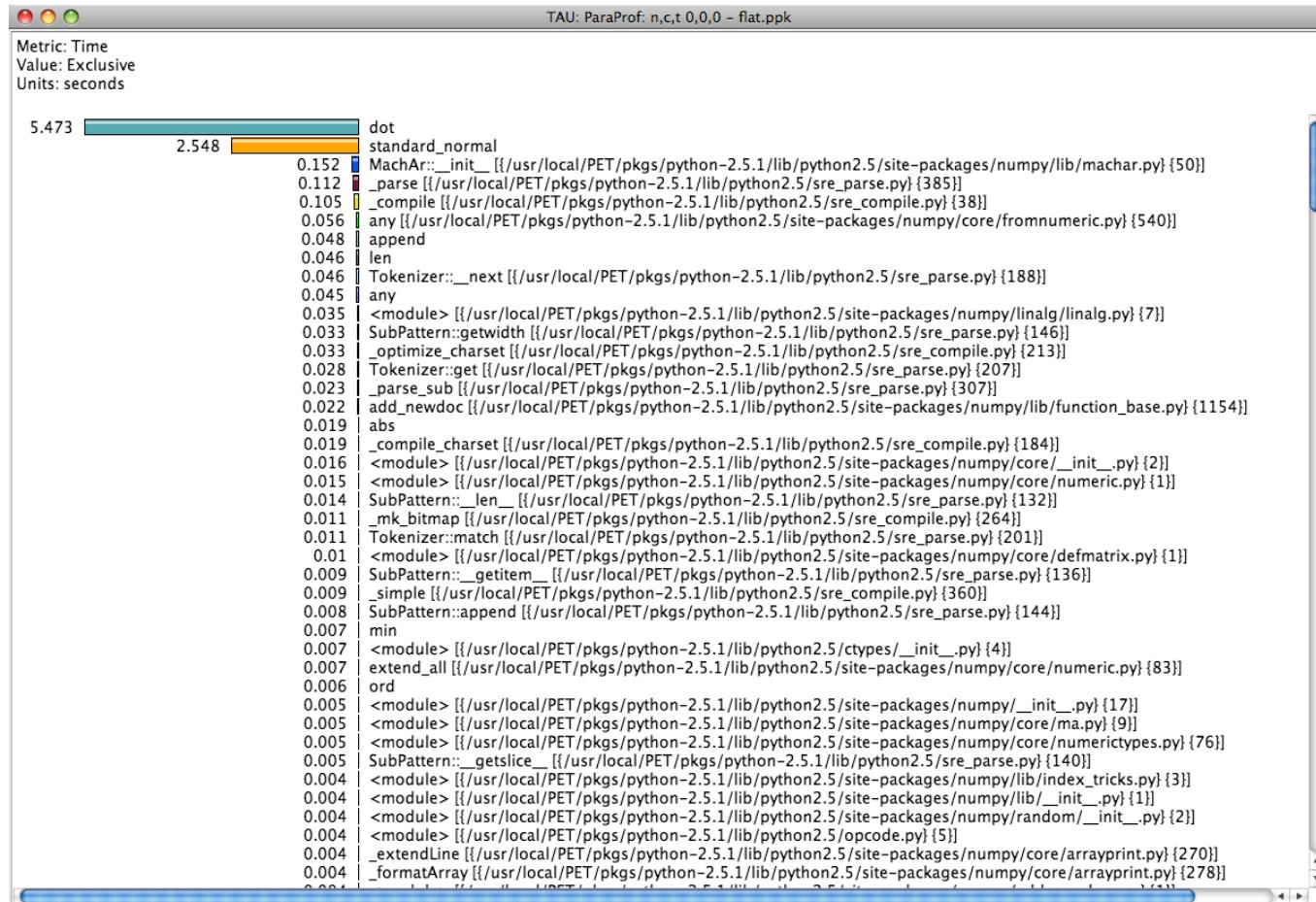


Detect Memory Leaks

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
                  /lib/Makefile.tau-mpi-pdt  
  
% setenv TAU_OPTIONS '-optDetectMemoryLeaks -optVerbose'  
  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
  
% make F90=tau_f90.sh  
  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% setenv TAU_CALLPATH_DEPTH 100  
  
  
% qsub run.job  
  
% paraprof --pack app.ppk  
  Move the app.ppk file to your desktop.  
  
% paraprof app.ppk  
  (Windows -> Thread -> Context Event Window -> Select thread -> select...  
   expand tree)  
  (Windows -> Thread -> User Event Bar Chart -> right click LEAK  
   -> Show User Event Bar Chart)
```

Usage Scenarios: Instrument a Python program

- Goal: Generate a flat profile for a Python program



Usage Scenarios: Instrument a Python program

*Original
code:*

```
% cat foo.py
#!/usr/bin/env python
import numpy
ra=numpy.random
la=numpy.linalg

size=2000
a=ra.standard_normal((size,size))
b=ra.standard_normal((size,size))
c=la.linalg.dot(a,b)
print c
```

Create a wrapper:

```
% cat wrapper.py
#!/usr/bin/env python

# setenv PYTHONPATH $PET_HOME/pkgs/tau-2.17.3/ppc64/lib/bindings-gnu-python-pdt

import tau

def OurMain():
    import foo

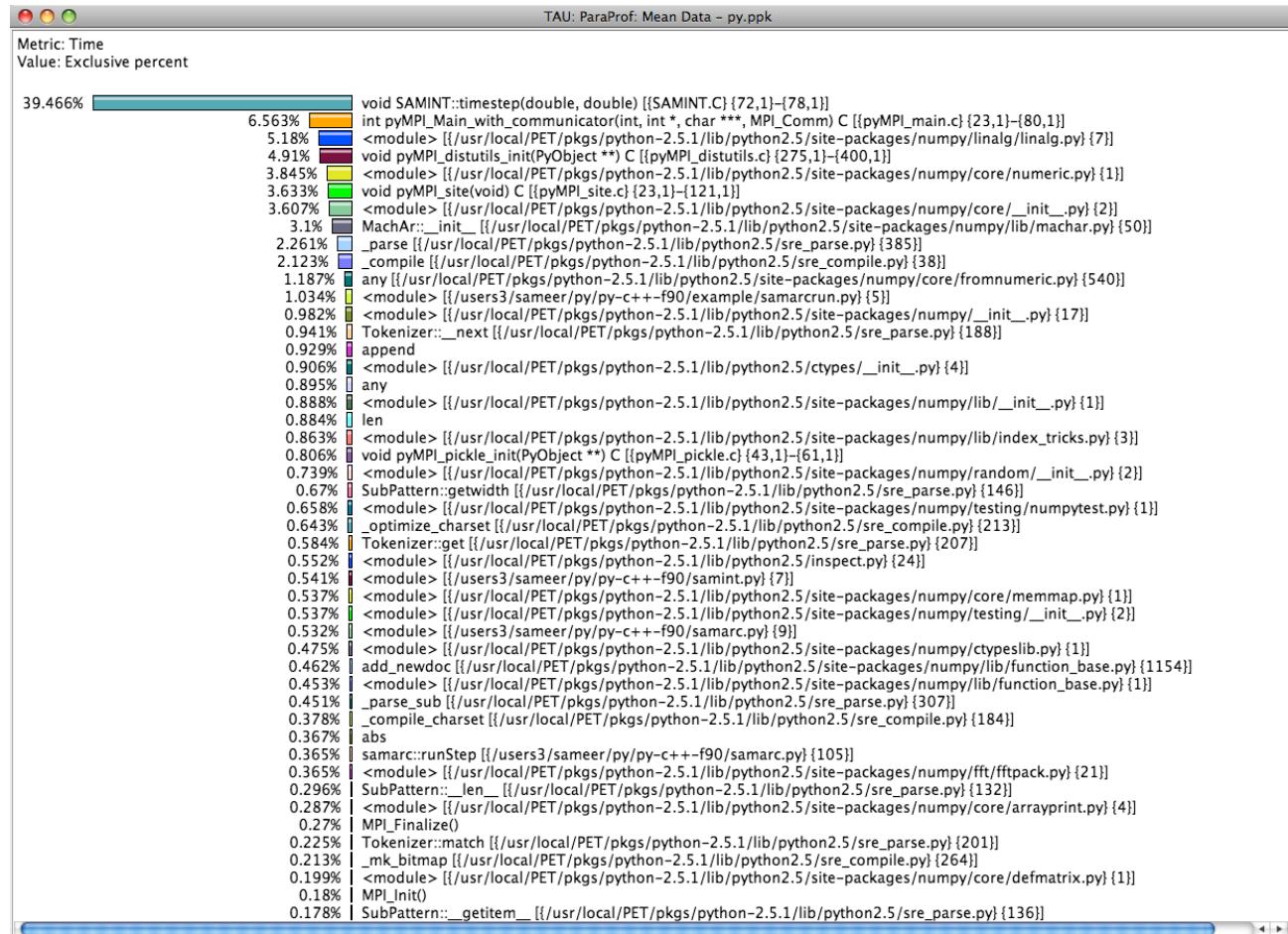
tau.run('OurMain()')
```

Generate a Python Profile

```
% setenv TAU_MAKEFILE /usr/global/tools/tau/ibm64  
                  /lib/Makefile.tau-python-pdt  
% set path=(/usr/global/tools/tau/ibm64/bin $path)  
% cat wrapper.py  
import tau  
def OurMain():  
    import foo  
    tau.run('OurMain()')  
Uninstrumented:  
% ./foo.py  
Instrumented:  
% setenv PYTHONPATH <taudir>/ibm64/<lib>/bindings-python-pdt  
(same options string as TAU_MAKEFILE)  
% setenv LD_LIBRARY_PATH <taudir>/x86_64/lib/bindings-python-pdt\  
$LD_LIBRARY_PATH  
% ./wrapper.py  
  
Wrapper invokes foo and generates performance data  
% pprof/paraprof
```

Usage Scenarios: Mixed Python+F90+C+pyMPI

- Goal: Generate multi-level instrumentation for Python+MPI+C+F90+C++ ...

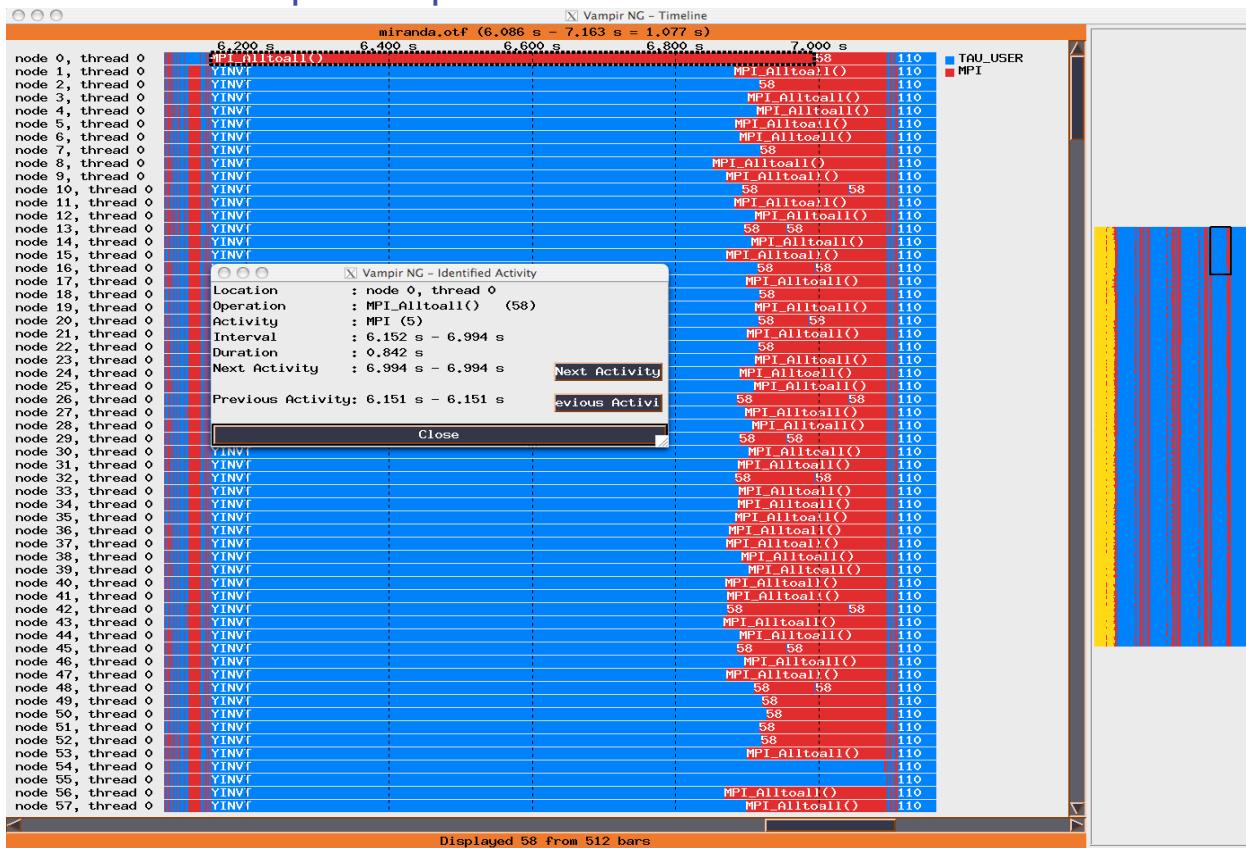


Generate a Multi-Language Profile w/ Python

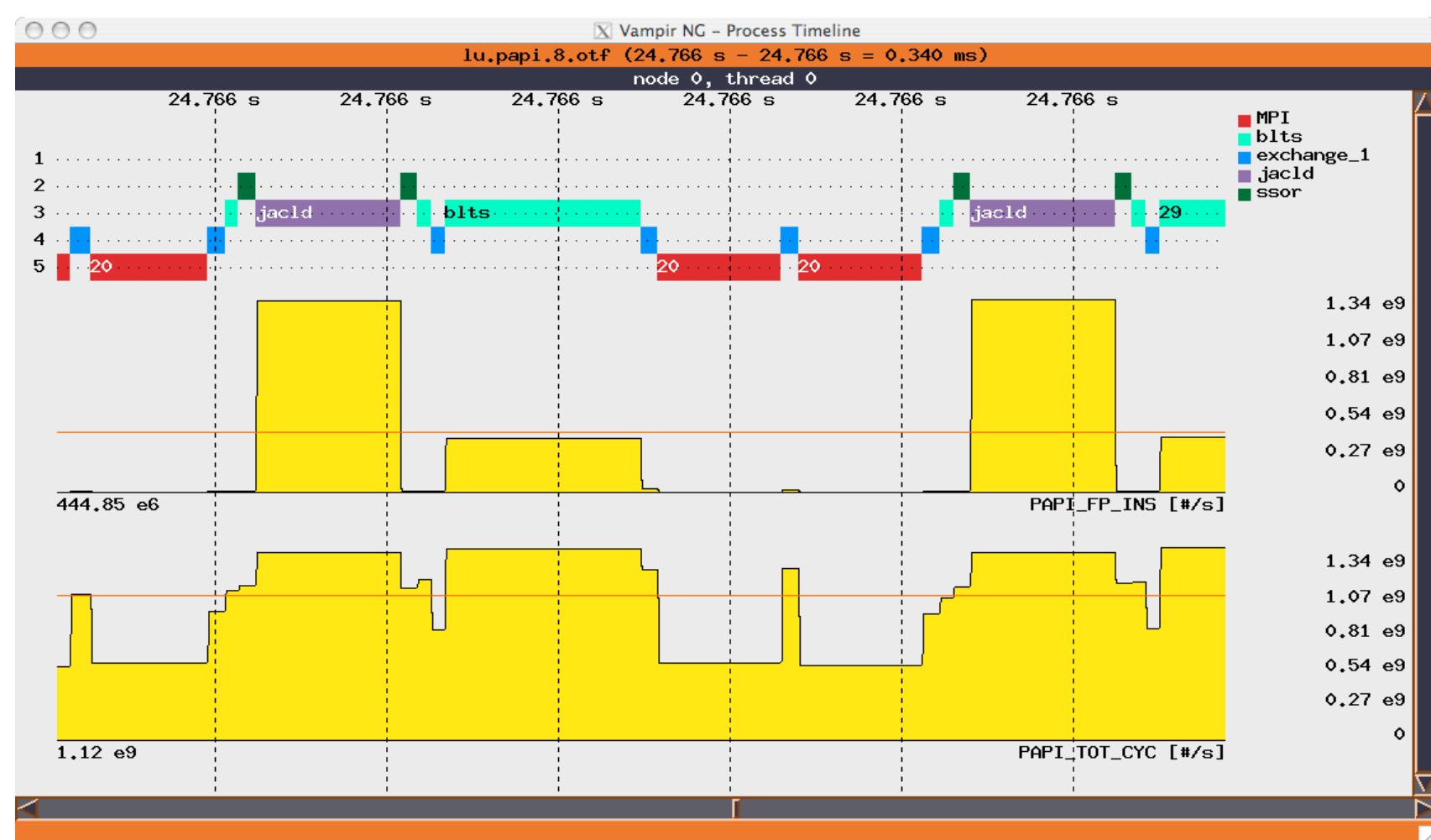
```
% setenv TAU_MAKEFILE /usr/global/tools/tau/ibm64  
                  /lib/Makefile.tau-python-mpi-pdt  
% set path=(/usr/global/tools/tau/ibm64/bin $path)  
% setenv TAU_OPTIONS '-optShared -optVerbose...'  
(Python needs shared object based TAU library)  
% make F90=tau_f90.sh CXX=tau_cxx.sh CC=tau_cc.sh  (build libs, pyMPI w/TAU)  
% cat wrapper.py  
import tau  
  
def OurMain():  
    import App  
    tau.run('OurMain()')  
Uninstrumented:  
% mpirun.lsf /usr/global/tools/.unsupported/pyMPI-2.5b0/bin/pyMPI ./App.py  
Instrumented:  
% setenv PYTHONPATH <taudir>/x86_64/<lib>/bindings-python-mpi-pdt  
(same options string as TAU_MAKEFILE)  
% setenv LD_LIBRARY_PATH <taudir>/x86_64/lib/bindings-python-mpi-pdt\:$LD_LIBRARY_PATH  
% mpirun -np 4 /usr/global/tools/.unsupported/pyMPI-2.5b0-TAU/bin/pyMPI  
./wrapper.py          (Instrumented pyMPI with wrapper.py)
```

Usage Scenarios: Generating a Trace File

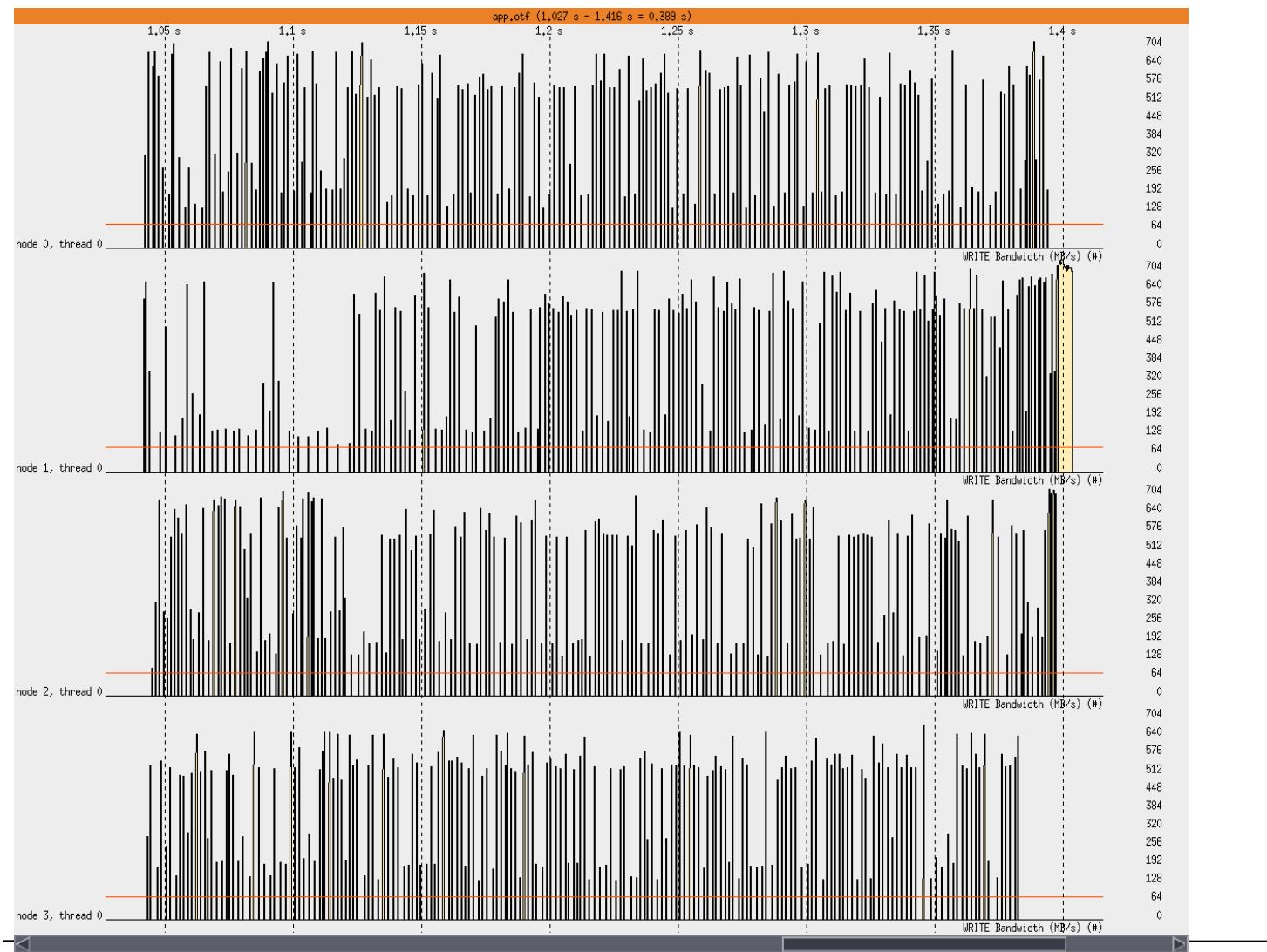
- Goal: Identify the temporal aspect of performance. What happens in my code at a given time? When?
- Event trace visualized in Vampir/Jumpshot



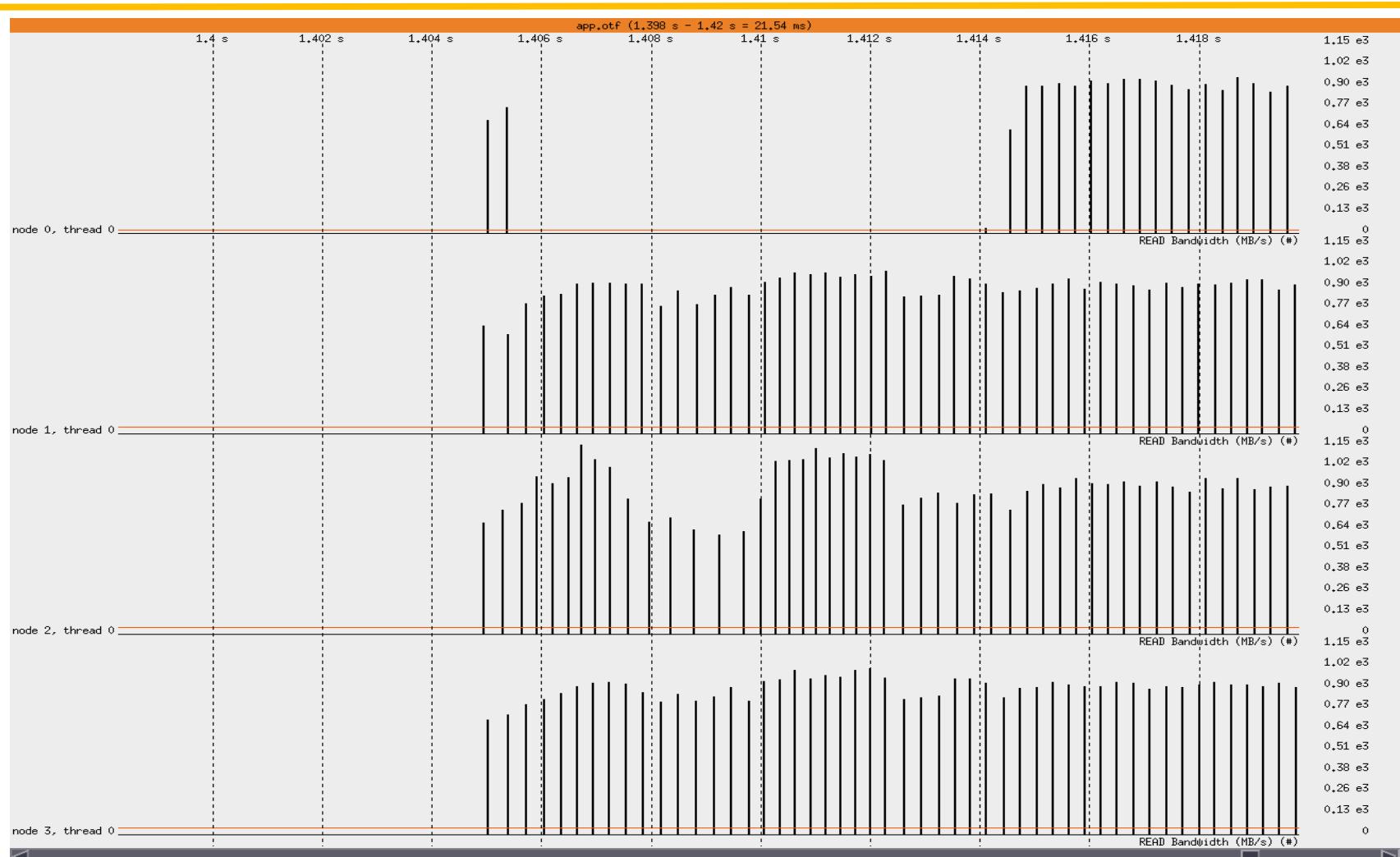
VNG Process Timeline with PAPI Counters



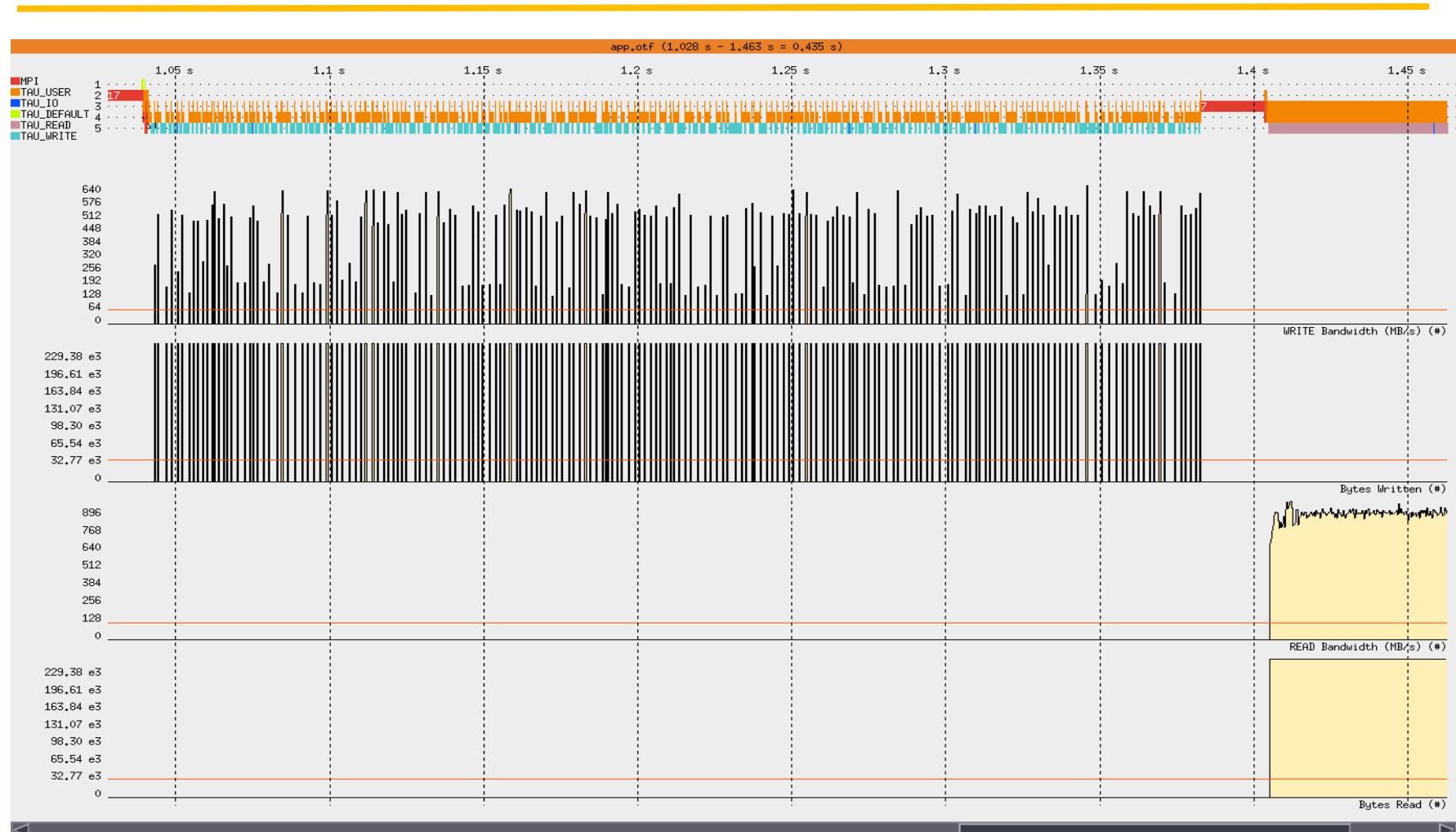
Vampir Counter Timeline Showing I/O BW



TAU: I/O Instrumentation for Read Bandwidth: Vampir



Vampir Process Timeline for Rank 0 (IOR, LLNL)

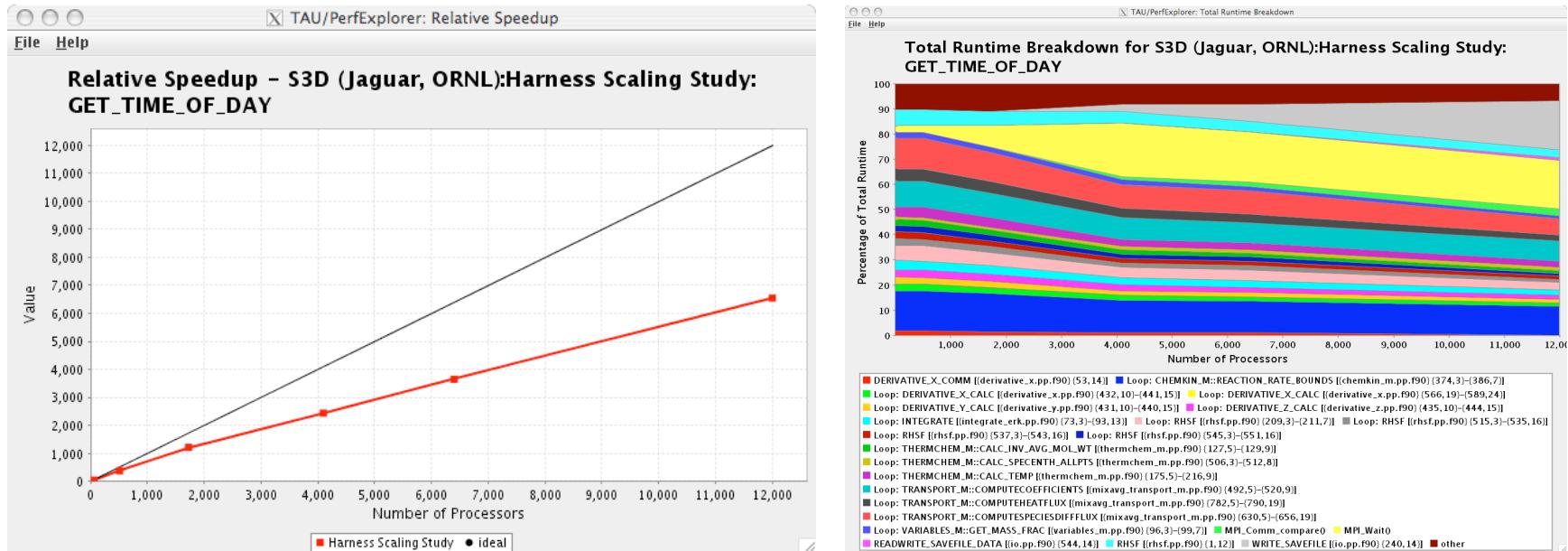


Generate a Trace File

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
                  /lib/Makefile.tau-mpi-pdt-trace  
or setenv TAU_TRACE 1 (in TAU v2.18.2+)  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
% make F90=tau_f90.sh  
(Or edit Makefile and change F90=tau_f90.sh)  
% qsub run.job  
% tau_treemerge.pl  
(merges binary traces to create tau.trc and tau.edf files)  
JUMPSHOT:  
% tau2slog2 tau.trc tau.edf -o app.slog2  
% jumpshot app.slog2  
    OR  
VAMPIR:  
% tau2otf tau.trc tau.edf app.otf -n 4 -z  
(4 streams, compressed output trace)  
% vampir app.otf  
(or vng client with vngd server).
```

Usage Scenarios: Evaluate Scalability

- Goal: How does my application scale? What bottlenecks occur at what core counts?
- Load profiles in PerfDMF database and examine with PerfExplorer

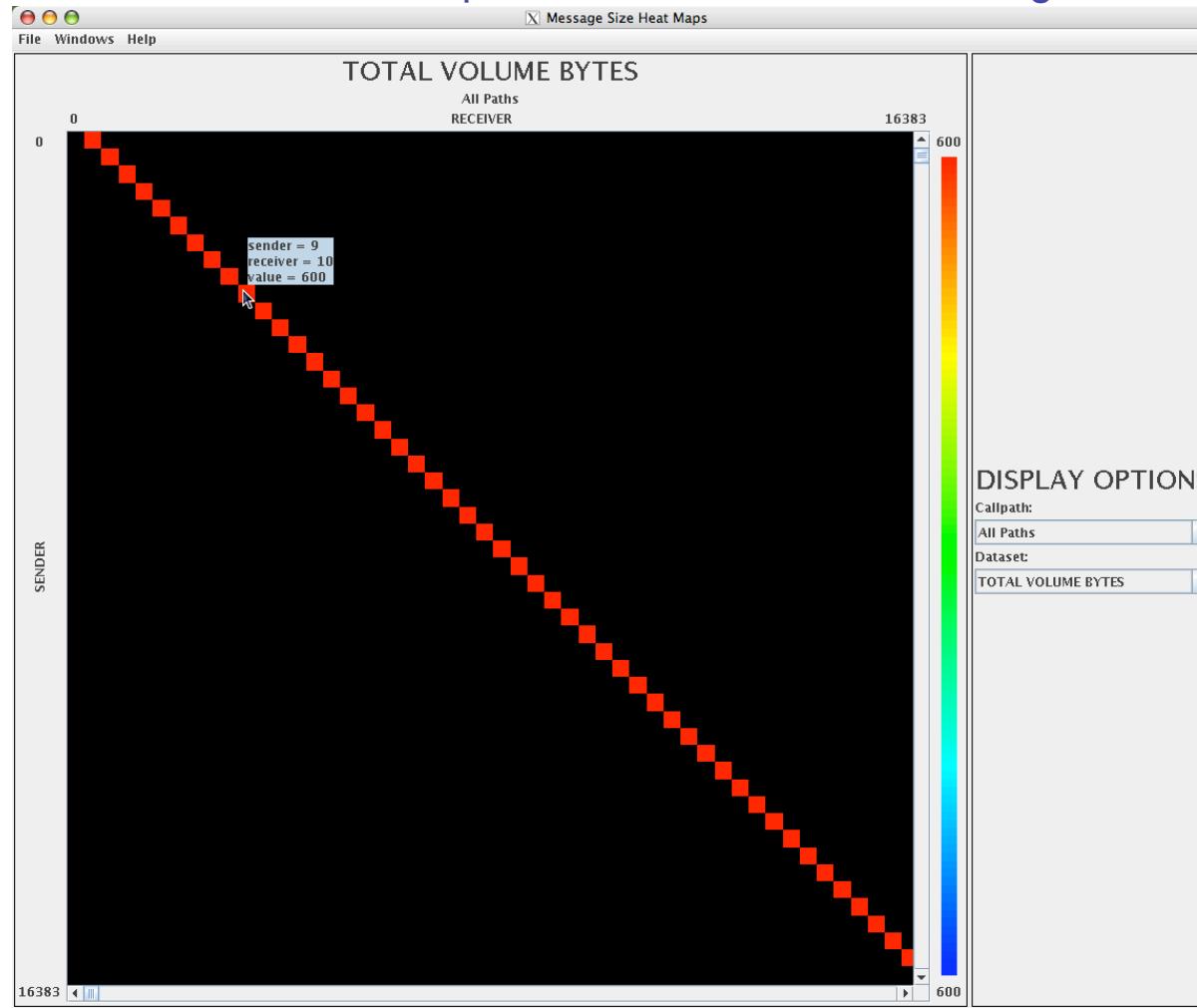


Evaluate Scalability using PerfExplorer Charts

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
          /lib/Makefile.tau-mpi-pdt  
  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
  
% make F90=tau_f90.sh  
  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% qsub run1p.job  
  
% paraprof --pack 1p.ppk  
  
% qsub run2p.job ...  
  
% paraprof --pack 2p.ppk ... and so on.  
  
On your client:  
  
% perfdmf_configure --create-default  
(Chooses derby, blank user/passwd, yes to save passwd, defaults)  
  
% perfexplorer_configure  
(Yes to load schema, defaults)  
  
% paraprof  
(load each trial: DB -> Add Trial -> Type (Paraprof Packed Profile) -> OK) OR use  
    perfdmf_loadtrial  
  
Then,  
  
% perfexplorer  
(Select experiment, Menu: Charts -> Speedup)
```

Communication Matrix Display

- Goal: What is the volume of inter-process communication? Along which calling path?



Evaluate Scalability using PerfExplorer Charts

```
% setenv TAU_MAKEFILE /soft/apps/tau/tau-2.18.2/bgp  
          /lib/Makefile.tau-mpi-pdt  
  
% set path=(/soft/apps/tau/tau-2.18.2/ppc64/bin $path)  
  
% make F90=tau_f90.sh  
  
(Or edit Makefile and change F90=tau_f90.sh)  
  
% setenv TAU_COMM_MATRIX 1  
  
  
% qsub run.job (setting the environment variables)  
  
  
% paraprof  
  
(Windows -> Communication Matrix)
```

Labs

- Add one of
source /soft/apps/tau/src/tau.bashrc
or
source /soft/apps/tau/src/tau.cshrc
to the end of your **.login** file (for bash or csh/tcsh users respectively)
These files contain ANL specific location information.

Support Acknowledgements

- Department of Energy (DOE)
 - Office of Science
 - MICS, Argonne National Lab
 - ASC/NNSA
 - University of Utah ASC/NNSA Level 1
 - ASC/NNSA, LLNL
- Department of Defense (DoD)
 - HPC Modernization Office (HPCMO)
- NSF SDCI
- Research Centre Juelich
- LBL, ORNL, ANL, LANL, LLNL
- TU Dresden
- University of Oregon

ParaTools

